

# **HARDNESS AND TRACTABILITY FOR STRUCTURED NUMERICAL PROBLEMS**

A Dissertation  
Presented to  
The Academic Faculty

By

Peng Zhang

In Partial Fulfillment  
Of the Requirements for the Degree  
Doctor of Philosophy in  
Computer Science

Georgia Institute of Technology

December 2018

Copyright © Peng Zhang 2018

# **HARDNESS AND TRACTABILITY FOR STRUCTURED NUMERICAL PROBLEMS**

Approved by:

Dr. Richard Peng, Advisor  
School of Computer Science  
*Georgia Institute of Technology*

Dr. Edmond Chow  
School of Computational Science  
and Engineering  
*Georgia Institute of Technology*

Dr. Petros Drineas  
Department of Computer Science  
*Purdue University*

Dr. Milena Mihail  
School of Computer Science  
*Georgia Institute of Technology*

Dr. Santosh Vempala  
School of Computer Science  
*Georgia Institute of Technology*

Date Approved: August 20 , 2018

To my parents.

## ACKNOWLEDGEMENTS

First of all, I am grateful to my advisor, Richard Peng, for his support and patience during my PhD study at Georgia Tech. Richard introduced me to these interesting theories and problems in linear systems and linear programs. He also provided me many invaluable advices in both research and career.

Second, I would like to thank my thesis committee members: Edmond Chow, Petros Drineas, Milena Mihail, and Santosh Vempala. I also thank Eric Vigoda for being my qualifying exam committee member.

I am lucky to work with my co-author, Rasmus Kyng. Most of the works in this thesis were collaborated with him, from which I learnt a lot. I would like to thank Di Wang, for working together on positive LPs. I also enjoyed talking research with Anup Rao, who answered my various silly questions during my first year at GT.

I would like to thank Dan Spielman, for his kind help during my visit at Yale. Working with Dan is amazing, full of fun and surprises. Then, I would like to thank Gary Miller, for the wonderful summer at CMU. I really enjoyed talking about geometry and distributions. I also thank Yi Wu, for guiding me into numerical linear algebra and teaching me many basic things.

I thank Ahmed Sameh and Luca Dieci for teaching me linear algebra. I am also glad of being in classes taught by Santosh Vempala, Santanu Dey, and Robin Thomas.

My two and half years at Georgia Tech were great with friends: Zongchen Chen, Tung Mai, Samira Samadi, David Durfee, Saurabh Sawlani, Matthew Fahrbach, Sam Petti, Kevin Lai, Chunxing Yin, Zhuangdi Xu, Tianxin Tang, Shan Chen, and too many to mention. Besides, I would like to thank H. Venkateswaran and Dani Denton for their help on my PhD program. I also thank Richard Lipton for his books at theory lab.

Outside GT, I owe my thanks to people at Purdue: Elena Grigorescu, Greg Frederickson, Mikhail Atallah, and David Gleich, for their generous help during my two years there.

I also thank my colleagues and friends in theory group: Akash Kumar, Mingyuan Wang, Abram Magner, Abhiram Natarajan, G.V., Samson Zhou, and Young-San Lin.

I would like to thank people in China: Wei Chen, Xiaoming Sun, and Guochuan Zhang, for introducing me to theoretical computer science, and for their advices on research and career.

My special thanks go to Zichao Yang and Li Xiong, for talking and complaining during the five years in the US.

Last but not the least, I would like to thank my parents, for their lifetime support and love.

## TABLE OF CONTENTS

<b>Acknowledgments</b> . . . . .	iv
<b>List of Figures</b> . . . . .	x
<b>Chapter 1: Introduction</b> . . . . .	1
1.1 Our Results . . . . .	7
1.2 Organization of This Dissertation . . . . .	12
<b>Chapter 2: Preliminaries</b> . . . . .	13
2.1 Vectors and Matrices . . . . .	13
2.2 Approximately Solving A Linear System . . . . .	15
2.3 Graph-Structured Block Matrices . . . . .	18
2.4 Geometric Structures for 3-D Trusses . . . . .	21
2.5 Techniques for Solving Linear Systems . . . . .	23
<b>Chapter 3: Hardness Results of Structured Linear Systems</b> . . . . .	26
3.1 Measuring the difficulty of solving a linear system . . . . .	26
3.2 Main Results and Reduction Outline . . . . .	28
3.3 Reducing Zero-Sum Power Two Linear Systems to Two-Commodity Linear Systems . . . . .	30
3.3.1 Approximate solvers . . . . .	43

3.3.2	Bounding Condition Number of the New Matrix . . . . .	46
3.3.3	Putting it All Together . . . . .	58
3.4	$\mathcal{MC}_2$ Efficiently Reducible to $\mathcal{MC}_2^{>0}$ . . . . .	59
3.4.1	Construction . . . . .	59
3.4.2	Bounding Condition Number of the New System in $\mathcal{MC}_2^{>0}$ . . . . .	66
3.4.3	Putting it All Together . . . . .	67
3.5	Rounding and Scaling Weights to Integers . . . . .	68
3.6	$\mathcal{G}$ Efficiently Reducible to $\mathcal{G}_{z,2}$ . . . . .	74
3.6.1	$\mathcal{G} \leq_f \mathcal{G}_z$ . . . . .	75
3.6.2	$\mathcal{G}_z \leq_f \mathcal{G}_{z,2}$ . . . . .	82
3.7	2D Trusses . . . . .	93
3.8	Connections with Interior Point Methods . . . . .	99
3.8.1	Brief Overview of Interior Point Methods . . . . .	99
3.8.2	2-Commodity Flow . . . . .	100
3.8.3	Isotropic Total Variation Minimization . . . . .	111
<b>Chapter 4: Hardness Results of Packing LPs . . . . .</b>		<b>116</b>
4.1	Reducing a Linear System to a Packing LP . . . . .	117
4.1.1	Dense Reduction $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$ . . . . .	120
4.1.2	Sparse Reduction $\text{PLP}_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$ . . . . .	125
4.1.3	Reducing $\text{poly}(1/\epsilon)$ -Dependence to $\log(1/\epsilon)$ -Dependence . . . . .	129
<b>Chapter 5: 3-D Truss Linear System Solver . . . . .</b>		<b>131</b>
5.1	Main Algorithm and Proofs of Main Results . . . . .	131

5.2	Bounding the Smallest Nonzero Eigenvalue of a edge-simple Truss . . . . .	137
5.2.1	Main Ideas . . . . .	137
5.2.2	Centering a Vector (Proof of Lemma 5.2.1) . . . . .	140
5.2.3	Existence of Good Centering (Proof of Lemma 5.2.3) . . . . .	144
5.3	Proof of Path Lemma . . . . .	150
5.3.1	Relabeling Tetrahedrons and Vertices . . . . .	150
5.3.2	Distance between Local Minimizers and the Centered Vectors . . .	151
5.3.3	Bounding the Norm of $\bar{q}_i - \bar{q}_{\sigma_i(j)}$ in terms of $d_i$ 's . . . . .	157
5.3.4	Bounding the Quadratic Form $(\bar{q})^\top M \bar{q}$ . . . . .	162
5.4	Computing a $(B, r)$ -Hollowing of a Convex Edge-simple Truss . . . . .	163
5.4.1	Bounding the Size of $\mathcal{H}$ . . . . .	164
5.4.2	Bounding the Number of Tetrahedrons in $\mathcal{H}$ Intersecting with a Plane	167
5.4.3	Bounding the Relative Condition Number of $\mathcal{H}$ . . . . .	171
5.5	Nested Dissection . . . . .	172
5.5.1	Proof of Lemma 2.5.2 . . . . .	172
5.5.2	Proof of Lemma 5.1.5 . . . . .	177
5.6	Running Times In Terms of Fast Matrix Multiplication . . . . .	181
<b>Chapter 6: Parallel Algorithm for Mixed Packing and Covering LPs . . . . .</b>		<b>185</b>
6.1	Converting an Optimization Problem to Feasibility Problems . . . . .	185
6.2	Algorithm for the Feasibility Problem . . . . .	186
6.2.1	Potential Function . . . . .	186
6.2.2	Algorithm . . . . .	187



6.3	Analysis . . . . .	188
6.3.1	Proof of Correctness . . . . .	188
6.3.2	Analysis of Convergence . . . . .	194
6.3.3	Pure Packing and Pure Covering LPs . . . . .	200
<b>Appendix A: Proofs of Linear Algebra Facts . . . . .</b>		<b>202</b>
A.1	Schur Complements . . . . .	202
A.2	Solutions of Linear Systems . . . . .	203
<b>References . . . . .</b>		<b>212</b>

## LIST OF FIGURES

3.1	The relation between vectors $\tilde{\mathbf{y}}, -\alpha \mathbf{1}$ and $\tilde{\mathbf{y}} - \alpha \mathbf{1}$ . . . . .	81
3.2	Geometric realization of the multicommodity flow gadget generated by as a truss matrix . . . . .	95
5.1	An example of the relabeling of tetrahedrons and vertices: $t_1 = \{-2, -1, 0, 1\}, t_2 = \{-2, -1, 1, 2\}, t_3 = \{-2, 0, 1, 3\}, t_4 = \{-1, 0, 1, 4\}$ and $\sigma_1 = \{-2, -1, 0\}, \sigma_2 = \{-2, -1, 1\}, \sigma_3 = \{-2, 1, 0\}, \sigma_4 = \{0, -1, 1\}$ . . . . .	151
5.2	Tetrahedrons $t_{i_1}$ and $t_{i'}$ . . . . .	158
5.3	If $u_1$ is a child of $u$ in $T_{\text{rec}}^{(i,j)}$ , then the depth of $t_{i_{u_1}}$ is smaller than the depth of $t_{i_u}$ in $T_{\text{BFS}}$ . . . . .	159
5.4	An example of $\phi(f)$ . Face $f$ has five vertices $v_1, \dots, v_5$ . $\phi(f)$ is the surface of $B'_1$ enclosed by $(u_1, \dots, u_5, u_1)$ . . . . .	166
5.5	Plane $P$ intersects face $B_{x1}$ with line $(K, L)$ . We draw a line going through point $K$ and parallel to the $x$ -axis, which intersects face $B_{x2}$ at point $M$ , similarly we draw a line going through point $L$ and parallel to the $x$ -axis, which intersects face $B_{x2}$ at point $N$ . . . . .	168
5.6	Cut the shape along the plane $P$ , and shift the left part along the $x$ -axis and glue the two faces $KFJGL$ and $MEIHN$ . . . . .	169

## SUMMARY

We study structured linear systems and structured linear programs from both algorithm and complexity perspectives. These structured problems commonly arise in combinatorial optimization, machine learning, and operation research. Many of them can be solved significantly faster than general linear systems or linear programs by utilizing additional structures.

First, we consider linear systems in matrices which form a slightly larger class of graph Laplacians. Nearly-linear time algorithms have been designed for linear systems in graph Laplacians (Spielman-Teng, 2004), and sequentially for its generalizations, e.g., connection Laplacians (Kyng et al, 2016), and directed Laplacians (Cohen et al, 2017, 2018). In contrast to these positive results, we establish lower bounds for approximately solving linear systems in: multi-commodity Laplacian matrices, truss stiffness matrices, total variation matrices, all of which fall into the class of generalized graph Laplacians. Specifically, approximately solving linear systems in these matrices is equally hard as solving arbitrary linear systems. Our result answers a major question in this field (Spielman, 2016).

As a follow up, we prove hardness results for packing and covering LPs, which are special case of LPs with non-negative coefficients, constants and variables. We show that: any approximate solvers for packing and covering LPs whose run time has small dependence on error parameter would directly improve the best solvers for arbitrary linear systems. The latter is a long-standing open question in numerical linear algebra and scientific computing.

On the algorithmic side, we study linear systems in 3-D truss stiffness matrices with additional geometric structures. By combining Nested Dissection and support theory, which are both well-studied techniques for linear systems, we obtain an asymptotically faster algorithm. Moreover, we design a parallel algorithm for approximately solving mixed packing and covering LPs, which improves the algorithm by Young (2001).

# CHAPTER 1

## INTRODUCTION

A linear program (LP) is an optimization problem with linear constraints and linear objective function. In general, it can be written as:

$$\max_x \{ \mathbf{c}^\top \mathbf{x} : \mathbf{A}\mathbf{x} \leq \mathbf{b} \}, \quad (1.1)$$

where  $\mathbf{A}$  is a given matrix,  $\mathbf{b}$ ,  $\mathbf{c}$  are given vectors, and  $\mathbf{x}$  is an unknown vector. There are many other equivalent forms of linear programs. A linear system is a system of linear equations. Its input consists of a coefficient matrix  $\mathbf{A}$ , a column vector  $\mathbf{b}$ , and its output is a vector  $\mathbf{x}$  such that<sup>1</sup>

$$\mathbf{A}\mathbf{x} = \mathbf{b}.$$

In practice, we allow small error  $\epsilon > 0$ , and then the goal becomes to compute an approximate vector  $\tilde{\mathbf{x}}$  such that

$$\|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b}\|_2 \leq \epsilon \|\mathbf{b}\|_2.$$

There are several different definitions of approximately solving an LP, which we will not discuss in this thesis. We suggest readers to refer to [1, 2].

Linear programs and linear systems are central tools in computer science, combinatorial optimization, statistics, economics, and engineering. A variety of combinatorial problems, such as network flow problems, multicommodity flow problems, scheduling problems, etc., can be formulated as linear programs. One widely studied problem for modeling in statistics is linear regression problem. Given a (tall) matrix  $\mathbf{A}$ , a column vector  $\mathbf{b}$ , it seeks a vector  $\mathbf{x}$  to minimize  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_p$  for some norm  $p$ . When  $p = 1$  or  $\infty$ , this problem is

---

<sup>1</sup>Here we assume  $\mathbf{A}$  is invertible, we will discuss more general cases in Section 2.2.

equivalent to solving a linear program; when  $p = 2$ , this problem is equivalent to solving a linear system, which is usually referred to as the normal equation. Besides, both linear programs and linear systems are frequently used as sub-routines for other optimization problems.

Linear programs and linear systems are related via Interior-Point Method, which solves a linear program in provably polynomial time [3]. Interior-point method reduces solving a linear program of the form (1.1) to solving a sequence of  $\tilde{O}(\sqrt{\text{rank}(\mathbf{A})}L)^2$  linear systems, where  $L$  is the bit complexity of the input [4]. Most combinatorial problems have polynomially bounded input numbers, and thus  $L = O(\log n)$  is ignorable. Thus, any improvement on solving linear systems would imply a direct improvement on solving linear programs.

Interior-point method has advantages over other LP algorithms either in practice or in theory. Ellipsoid method is the first LP algorithm which was proved to have polynomial running time [5]. However, it has high cost in practice. The simplex algorithm runs well in practice, but there exists an example showing that the simplex algorithm needs exponential running time in the worst case [6]. About twenty years ago, it was shown that the simplex algorithm has smoothed polynomial time [7], which explains its good performance in practice.

In general, algorithms for solving linear programs or linear systems have high complexity in the worst case<sup>3</sup>. Asymptotically faster algorithms (even nearly-linear time algorithms) exist when inputs have some additional structures.

In this thesis, we study structured linear systems and linear programs from both algorithm and complexity perspectives. We consider linear systems in matrices which form a slightly larger class of graph Laplacians, referred to as Graph-Structured Block Matrices (GSBMs) [8]. In contrast to the existing nearly-linear time solvers for graph Laplacians [9] and its generalizations [10, 11, 12], we prove hardness results for GSBMs: approximately

---

<sup>2</sup>We use  $\tilde{O}(\cdot)$  to hide log factors.

<sup>3</sup>It is open whether linear programs are solvable in strongly-polynomial time.

solving linear systems in GSBMs is equally hard as approximately solving arbitrary linear systems. Building upon linear system based hardness assumptions, we then establish conditional lower bounds for packing and covering LPs, which are a special case of LPs with non-negative coefficients, constants and variables.

On the algorithmic side, we obtain an asymptotically faster solver for linear systems in 3-D trusses with additional geometric structures. General truss matrices are GSBMs. Moreover, we design a parallel algorithm for approximately solving mixed packing and covering LPs, which improves the algorithm by Young [13].

**Solving (Arbitrary) Linear Systems.** Given a linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  where  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , one can directly solve it by Gaussian elimination, which has running time  $O(n^3)$ . Thanks to the fast matrix multiplication since Strassen [14], one can speed up solving the linear system by inverting the coefficient matrix  $\mathbf{A}$  and then multiplying  $\mathbf{A}^{-1}\mathbf{b}$ . It was shown that one can invert an  $n \times n$  nonsingular matrix by recursively applying matrix multiplication, which takes  $O(n^\omega)$  time [15]. Here,  $\omega < 2.3728639$  [14, 16, 17] is the matrix multiplication constant. Using a similar trick, one can multiply two matrices by computing matrix inversions  $\log$  times. Thus, matrix multiplication and matrix inversion have same complexity up to a  $\log$  factor.

Note that  $\mathbf{A}^{-1}$  can be dense, even if  $\mathbf{A}$  itself is sparse. Thus, the matrix multiplication constant  $\omega \geq 2$ . Raz showed that multiplying two  $n \times n$  matrices needs time  $\Omega(n^2 \log n)$  [18]. There is no other non-trivial lower bound for  $\omega$ .

This simple algorithm by fast matrix inversion is the best known algorithm for solving an arbitrary linear system. It is not known whether one can solve an arbitrary linear system in  $n \times n$  dimensions in time  $o(n^\omega)$ .

When a linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  is sparse (e.g.,  $\mathbf{A}$  has  $n \log n$  non-zeros), it is desirable to solve this linear system in time slightly larger than  $\text{nnz}(\mathbf{A})$ , the number of non-zeros of  $\mathbf{A}$ <sup>4</sup>. Iterative methods aim at this goal, by avoiding matrix inversion. The operations of

---

<sup>4</sup>One needs to read all entries of  $\mathbf{A}$  to solve a linear system with small error.

each iteration only include matrix-vector multiplication and vector-vector operations. The best known iterative method for solving an arbitrary linear system is Conjugate Gradient Descent or its variants [19] (e.g., Chebyshev Iterations, or an accelerated first order method such as Nesterov’s algorithm). These methods solve a linear system whose coefficient matrix  $\mathbf{A}$  is symmetric and PSD in time  $O(\text{nnz}(\mathbf{A})\sqrt{\kappa(\mathbf{A})}\log(1/\epsilon))$ , where  $\kappa(\mathbf{A})$  is the condition number of  $\mathbf{A}$ .

Conjugate gradient can also be shown in the RealRAM model to converge to an exact solution in  $n$  iterations, giving a running time of order  $n \cdot \text{nnz}(\mathbf{A})$ . However, this is highly misleading, because achieving this type of behavior in a fixed point arithmetic model requires about  $n$  bits per number, as opposed to  $\text{poly} \log(n/\epsilon)$  bits to achieve the condition number dependent behavior [20].

**Graph Laplacians, its Generalizations, and Nearly-Linear Time Solvers.** Graph Laplacians are a class of well-studied structured matrices. Linear systems in graph Laplacians can be solved in nearly-linear time.

An  $n \times n$  Laplacian matrix can be viewed as an undirected graph over  $n$  vertices. Each vertex of the graph corresponds to a single row and a single column of the matrix. Each edge corresponds to an  $n$ -dimensional vector whose non-zeros are only at the position of its two endpoints. Taking the outer product of the edge vector, we get a rank-one matrix. The Laplacian matrix is then defined as a sum over these edge rank-one matrices.

Laplacian linear systems arise from using interior-point methods for combinatorial problems on graphs. For example, max flow in directed unweighted graphs [21, 22], negative weight shortest paths and max weight matchings [23], min cost flows and lossy generalized flows [24, 4], etc. Using fast Laplacian solvers to speeding up graph algorithms was formulated as the Laplacian Paradigm in [25].

Spielman and Teng designed the first *nearly-linear time* solver for Laplacian linear systems [9]. Specifically, their algorithms solves a linear system in a Laplacian matrix with

$m$  nonzeros up to accuracy  $\epsilon > 0$  in time  $O(m \log^{O(1)} n \log(1/\epsilon))$ . The exponent constant of  $\log n$  was then improved by a sequence of following works (see [26, 27, 28, 29, 30, 31], etc.).

A natural generalization of graph Laplacians is to allow *multiple* labels for each vertex. These matrices arise from using interior-point methods solving more generalized graph problems, e.g., multi-commodity flow problems [21, 22], and total variation minimization problem in images [32]. In addition, these matrices have been studied in Markov random fields [33], image processing [34, 35, 36], Euclidean embedding of graphs [37], data processing for cryo-electron microscopy [38, 39, 40], phase retrieval [41, 42], and linear elasticity [43, 44].

Connection Laplacians are one example, for which nearly-linear time solvers exist [10]. These matrices have applications in image and signal processing [45, 46]. Other examples include Truss stiffness matrices [43, 44], multi-commodity Laplacians [21, 22], and total variation minimization matrices [32]. However, it is not known whether one can solve a linear system in one of these matrices in nearly-linear time.

**Geometric Structures and Nested Dissection.** Another widely studied class of structured linear systems, for which asymptotically faster algorithms exist, is linear systems whose coefficient matrices' non-zeros have geometrically embedded structures. Given an  $n \times n$  matrix  $\mathbf{A}$ , its non-zero structures can be viewed as a graph over  $n$  vertices: There is an edge between vertices  $i$  and  $j$  iff the  $(i, j)$ th entry or the  $(j, i)$ th entry of  $\mathbf{A}$  is non-zero. We call this graph as  $G_A$ .

Applying Gaussian elimination to  $\mathbf{A}\mathbf{x} = \mathbf{b}$  can be viewed as operations on  $G_A$ : Eliminate the  $i$ th row and the  $i$ th column of  $\mathbf{A}$  corresponds to deleting the vertex  $i$  of  $G_A$  and putting a clique over the neighbors of  $i$ . This process may introduce new non-zeros in the matrix (correspondingly, new edges on the graph), which we refer to as fill-ins. The size of fill-ins directly determines the running time of Gaussian elimination. Even if  $\mathbf{A}$  itself



is sparse, the intermediate matrix produced by partial Gaussian elimination can be dense. The key to speeding up Gaussian elimination is to finding an elimination ordering which produces a small number of fill-ins.

*Nested Dissection* is used to compute an ordering for sparse Gaussian elimination when matrices  $\mathbf{A}$  with additional geometric structures, for example,  $G_A$  is a planar graph, a well-shaped 3-D tetrahedral meshes, or a graph without some minors [47, 48, 49, 50]. These graphs have small separators, whose removal divides a graph into balanced pieces. Nested dissection recursively computes small separators of  $G_A$ , and then eliminate vertices in a reverse order. Using this idea, one can solve a linear system in a 2-D planar in time  $O(n^{1.5})$  [48], and a linear system in a well-shaped 3-D tetrahedral mesh in time  $O(n^2)$  [49].

**Packing and Covering LPs.** Packing and covering LPs are LPs formulated with non-negative coefficients, constants and variables. These LPs commonly arise from theoretical computer science, operation research, and optimization, etc. Moreover, they can be approximately solved in nearly linear time.

Formally, a packing LP has its generic form

$$\max_{\mathbf{x} \geq \mathbf{0}} \{ \mathbf{c}^\top \mathbf{x} : \mathbf{A} \mathbf{x} \leq \mathbf{b} \},$$

where  $\mathbf{A}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$  have non-negative entries. A covering LP is of the form

$$\min_{\mathbf{x} \geq \mathbf{0}} \{ \mathbf{c}^\top \mathbf{x} : \mathbf{A} \mathbf{x} \geq \mathbf{b} \},$$

where  $\mathbf{A}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$  again are entrywise non-negative. One more general class is mixed packing and covering LPs, which is sometimes simply referred to as positive LPs. A positive LP in general has both packing constraints and covering constraints.

Note that a packing LP or a covering LP is always feasible, and has non-negative optimal value due to the all-zero solution. Besides, a packing LP and a covering LP form

a primal-dual pair. Let  $\text{OPT}$  denote the optimal value of this primal-dual pair. We say  $\mathbf{x}$  is a  $(1 - \epsilon)$ -approximation for the packing LP if  $\mathbf{x}$  is feasible and  $\mathbf{c}^\top \mathbf{x} \geq (1 - \epsilon)\text{OPT}$ . Similarly, we say  $\mathbf{x}$  is a  $(1 + \epsilon)$ -approximation for the covering LP if  $\mathbf{x}$  is feasible and  $\mathbf{c}^\top \mathbf{x} \leq (1 + \epsilon)\text{OPT}$ .

We are interested in *nearly-linear time* algorithms for approximately solving packing and covering LPs. Specifically, the running time is linear in the number of input non-zeros or dimensions, and polynomial in  $1/\epsilon$ . Note the running time is independent of input coefficients, and thus these algorithms are referred to as width-independent algorithms in the literature.

Although one can use ellipsoid methods or interior-point methods to get  $\log(1/\epsilon)$ -dependence [5, 3], these algorithms usually have large dependence on the dimension and the bit-complexity of inputs, which are not desirable in practice. The nearly-linear time solvers for packing and covering LPs are mainly based on first-order methods, e.g., gradient descent, multiplicative weights update [51, 13, 52, 53, 54, 55].

## 1.1 Our Results

**Hardness Results for Graph-Structured Linear Systems.** We consider three types of Graph-Structured Block Matrices (GSBM): Multi-commodity Laplacians, Truss Stiffness Matrices, and Total Variation Matrices. All these matrices are generalizations of graph Laplacians by allowing multiple labels for each vertex. We are interested in solving linear systems in GSBMs such that the running time has logarithmic dependence on condition number of coefficient matrices and accuracy. This is the kind of running time dependence established for graph Laplacians, Connection Laplacians, directed Laplacians, and planar 2-D truss stiffness matrices. The following result is joint with Kyng [56].

**Theorem 1.1.1** (Hardness for Graph-Structured Linear Systems). *We consider three types of Graph-Structured Block Matrices: 2-commodity Laplacians, 2-D Truss Stiffness Matrices, and 2-Total Variation Matrices. Suppose that for one or more of these classes, the*

linear system  $\mathbf{Ax} = \mathbf{b}$  in a matrix  $\mathbf{A}$  with  $s$  non-zeros can be solved in time  $\tilde{O}(s^a)$ , for some constant  $a \geq 1$ , with the running time having logarithmic dependence on condition number and accuracy. Then linear systems in all matrices with polynomially bounded integer entries and condition number can be solved to high accuracy in time  $\tilde{O}(s^a)$ , where again  $s$  is the number of non-zero entries of the matrix.

Particularly, if one can solve a linear system in a 2-commodity Laplacian, or a 2-D truss stiffness matrix, or a 2-total variation matrix in nearly linear time, as the existing solvers for graph Laplacians and connection Laplacians, then one can solve an arbitrary linear system with bounded entries and condition number in nearly linear time. However, the best known algorithm for solving an arbitrary linear system  $\mathbf{Ax} = \mathbf{b}$  of dimension  $n \times n$  runs in time  $O(\min\{n^\omega, \text{nnz}(\mathbf{A})\sqrt{\kappa(\mathbf{A})} \log(1/\epsilon)\})^5$ . This is asymptotically larger than any nearly linear function of  $\text{nnz}(\mathbf{A})$ , for  $\text{nnz}(\mathbf{A}) \ll n^2$ .

**Hardness Results for Packing LPs.** We have two hardness results for packing LPs. One concerns a dense packing LP, for which the running time is usually a function of dimensions. The other one is related to a sparse packing LP, for which the running time is a linear function of the number of non-zeros. Since packing LPs and covering LPs are primal and dual to each other, our results also hold for covering LPs. Both Theorem 1.1.2 and 1.1.3 are joint with Kyng, Peng and Di [57].

**Theorem 1.1.2.** *If one can  $(1 - \epsilon)$ -approximately solve a packing LP of dimension  $O(n) \times O(n)$  in time  $O(n^2/\epsilon^{0.0009})$ , then one can solve an arbitrary linear system of dimension  $O(n) \times O(n)$  asymptotically faster than multiplying two  $n \times n$  matrices.*

As mentioned earlier in Introduction, the best known algorithm for solving an arbitrary linear system of dimension  $n \times n$  runs in time  $O(n^\omega)$ , for  $\omega < 2.373$ . If  $\omega$  is strictly larger than 2 (e.g.,  $\omega > 2.00001$ ), then Theorem 1.1.2 would imply that “matrix inversion-by-matrix multiplication” is suboptimal to solving an arbitrary linear system.

---

<sup>5</sup>Or  $\text{nnz}(\mathbf{A})\kappa(\mathbf{A}) \log(1/\epsilon)$  if  $\mathbf{A}$  is not PSD.

Our second result is related to sparse systems.

**Theorem 1.1.3.** *If one can  $(1 - \epsilon)$ -approximately solve a packing LP with  $N$  non-zeros in time  $O(N/\epsilon^{0.22})$ , then one can solve an arbitrary linear system asymptotically faster than applying Conjugate Gradient.*

Conjugate Gradient (CG) is the best known algorithm for solving a sparse and well-conditioned linear system. The running of applying CG to a linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  with  $\mathbf{A}$  being symmetric and PSD is  $O(\text{nnz}(\mathbf{A})\sqrt{\kappa(\mathbf{A})}\log(1/\epsilon))$ , where  $\text{nnz}(\mathbf{A})$  is the number of non-zeros of  $\mathbf{A}$ ,  $\kappa(\mathbf{A})$  is the condition number, and  $\epsilon$  is the accuracy parameter. When  $\text{nnz}(\mathbf{A}) \approx n \log n$  and  $\kappa(\mathbf{A}) < n^{1.99}$ , CG is faster than directly inverting  $\mathbf{A}$ , even if the matrix multiplication constant  $\omega = 2$ .

Note for an arbitrary linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  satisfying the following:  $\dim(\mathbf{A}) = n \times n$ ,  $\text{nnz}(\mathbf{A}) = O(n \log n)$ , and  $\kappa(\mathbf{A}) = n^2$ , either fast matrix inversion or CG takes time  $\Omega(n^2)$ , given the fact that  $\omega \geq 2$ . It is not known whether one can approximately solve such a linear system in time  $o(n^2)$ . By Theorem 1.1.3, a fast packing LP solver would imply a positive answer to this question.

Trevisan and Xhafa [58] studied parallel complexity of packing LPs, and they showed that exactly solving a packing LP is P-complete. However, there is no known hardness result for approximately solving a packing LP. In [58], the authors gave an example showing that their reduction does not work for approximate solvers.

**Algorithms for Truss Linear Systems with Geometric Structures.** As stated in Theorem 1.1.1, solving an arbitrary 3-D truss linear system<sup>6</sup> is as hard as solving an arbitrary linear system. We consider 3-D trusses with additional geometric structures, which commonly arise from scientific computing.

A 3D truss is said to be *simple* iff it is a simplicial complex and every tetrahedron has bounded aspect ratio. We say a 3D truss is *edge-simple* iff its tetrahedral mesh is simple

---

<sup>6</sup>2-D truss is a special case of 3-D, by fixing one of the three dimensions to be 0.

and every tetrahedron has bounded edge lengths and stiffness coefficients (i.e., both are bounded above and below by constants). The following two results are from the paper [59] coauthored with Kyng, Peng and Schwieterman.

**Theorem 1.1.4.** *Consider an edge-simple 3-D truss  $\mathcal{T}$  over  $n$  vertices, which is formed by a union of  $k$  convex simplicial complexes with constant aspect ratio each. There is an algorithm which solves a linear system in the stiffness matrix of  $\mathcal{T}$  up to accuracy  $\epsilon > 0$  in time  $O(k^{1/3}n^{5/3} \log(1/\epsilon))^7$ .*

**Theorem 1.1.5.** *Consider an edge-simple 3-D truss  $\mathcal{T}$  over  $n$  vertices, which is formed by a union of  $k$  convex simplicial complexes with arbitrary aspect ratio each. There is an algorithm which solves a linear system in the stiffness matrix of  $\mathcal{T}$  up to accuracy  $\epsilon > 0$  in time  $O(k^{22/3}n^{11/6} \log(1/\epsilon))$ .*

Theorem 1.1.4 concerns a union of  $k$  convex truss with small aspect ratios, which improves the  $O(n^2)$  running time by nested dissection [49] for all  $k \ll n$ . Theorem 1.1.5 allows truss components to have arbitrary aspect ratios, which improves on nested dissection provided  $k \ll n^{1/44}$ .

Our requirements for trusses have more restrictions than the previous nested dissection based algorithms by Miller and Thurston [49], which only requires constant aspect ratios for individual tetrahedrons. Constant aspect ratios are crucial for nested dissection. Miller and Thurston provided a mesh of tetrahedrons with arbitrary aspect ratios, for which no sub-linear size separator exists and thus nested dissection cannot speed up Gaussian elimination. Our additional restrictions are mainly due to the need to derive spectral bounds for the preconditioner. Assumptions such as bounded edge lengths and stiffness coefficients (in addition to aspect ratios) are also present in previous work by Daich and Spielman [44], which forms the starting point for our key technical results on eigenvalues of 3-D simplicial complexes. Furthermore, the need to maintain null spaces as well as truss structures

---

<sup>7</sup>In both Theorem 1.1.4 and 1.1.5, we assume matrix multiplication constant  $\omega = 3$ , following the convention of nested dissection community. We will discuss the running times in term of the best known  $\omega < 2.3728639$  in Section 5.6.

in intermediate steps of Gaussian elimination places additional requirements on the mesh structure (mainly convexity).

We remark that the geometric assumptions of Theorem 1.1.5 are in many ways fairly weak. We need the individual truss tetrahedrons to have small aspect ratio, but each of the  $k$  truss components may overall have a wide range of shapes: It can form a ball, a pancake, or even a very long beam with arbitrarily large aspect ratio. The dependence on  $k$  is fairly bad and has not been carefully optimized, meaning currently that only about  $k \ll n^{1/44} \approx n^{0.0227}$  convex edge-simple trusses can be combined while still achieving a speed-up over nested dissection. Even so, this allows the construction of some shapes with genus up to  $n^{1/44}$ .

**Parallel Algorithm for Positive LPs.** We design the currently fast algorithm for approximately solving positive LPs in parallel. We use the following form of a positive LP:

$$\min_{\mathbf{x} \geq \mathbf{0}} \{ \lambda : \mathbf{P}\mathbf{x} \leq \lambda \mathbf{p}, \mathbf{C}\mathbf{x} \geq \mathbf{c} \},$$

where  $\mathbf{P}$ ,  $\mathbf{C}$ ,  $\mathbf{p}$ ,  $\mathbf{c}$  are all entrywise non-negative. A pair  $(\lambda, \mathbf{x})$  is a  $(1 + \epsilon)$ -approximation if it is a feasible solution and  $\lambda \leq (1 + \epsilon)\text{OPT}$ , where  $\text{OPT}$  is the optimal value of the above LP. The following result is a joint work with Mahoney, Rao and Wang [60].

**Theorem 1.1.6.** *Given a positive LP with  $N$  non-zero entries, there exists a parallel algorithm outputs a  $(1 + \epsilon)$ -approximate solution with running time  $\tilde{O}(1/\epsilon^3)$  and total work  $\tilde{O}(N/\epsilon^3)$ .*

This result improves Young’s algorithm in 2001 [13] by a factor  $1/\epsilon$ . Our algorithm follows the Lagrangian-relaxation approach, and can be viewed as a modification of Young’s algorithm. The bottleneck of each iteration is matrix-vector multiplication, which can be implemented parallelly in  $O(\log N)$  depth. Thus the running time is the number of iterations up to a log factor. The total work counts the number of total operations of the

algorithm. In general, the running time (or the number of iterations) is more interesting than the total work.

Recall that pure packing or covering LPs are a special case of positive LPs. It turns out that our algorithm approximately solves a pure packing or covering LP in  $\tilde{O}(1/\epsilon^2)$  iterations. This matches the best known results in [54, 61], while our algorithm is deterministic and has simpler analysis.

## **1.2 Organization of This Dissertation**

We give formal notations and definitions, which will be used throughout this thesis, in Chapter 2. We prove Theorem 1.1.1 in Chapter 3, Theorem 1.1.2 and 1.1.3 in Chapter 4, Theorem 1.1.4 and Theorem 1.1.5 in Chapter 5, and Theorem 1.1.6 in Chapter 6.

## CHAPTER 2

### PRELIMINARIES

#### 2.1 Vectors and Matrices

**Indexing.** We use subscripts to denote entries of a matrix or a vector. Given a vector  $\mathbf{x} \in \mathbb{R}^n$ , for  $1 \leq i < i+j \leq n$ , we denote  $\mathbf{x}_i$  the  $i$ th entry of  $\mathbf{x}$ , and we denote  $\mathbf{x}_{i:i+j}$  the subvector whose entries are  $\mathbf{x}_i, \mathbf{x}_{i+1}, \dots, \mathbf{x}_{i+j}$ . Given a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , let  $\mathbf{A}_i$  denote the  $i$ th row of  $\mathbf{A}$ , and let  $\mathbf{A}_{ij}$  denote the  $(i, j)$ th entry of  $\mathbf{A}$ . We use superscripts to index a sequence of matrices or vectors, e.g.,  $\mathbf{A}^1, \mathbf{A}^2, \dots$ , and  $\mathbf{x}^1, \mathbf{x}^2, \dots$ , except when some other meaning is clearly stated. In addition, we use  $\text{nnz}(\mathbf{A})$  to denote the number of nonzero entries of  $\mathbf{A}$ , and we use  $\dim(\mathbf{A})$  to denote the dimensions of  $\mathbf{A}$ .

**Spectral definitions.** A symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is *positive semi-definite* (PSD) iff  $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$  holds for every  $\mathbf{x} \in \mathbb{R}^n$ . For any matrix  $\mathbf{B} \in \mathbb{R}^{m \times n}$ , we denote  $\lambda_{\max}(\mathbf{B})$  the largest eigenvalue of  $\mathbf{B}$ , and  $\lambda_{\min}(\mathbf{B})$  the smallest *nonzero* eigenvalue of  $\mathbf{B}$ . Let  $\sigma_{\max}(\mathbf{B})$  be the largest singular value of  $\mathbf{B}$ , and  $\sigma_{\min}(\mathbf{B})$  be the smallest *nonzero* singular value of  $\mathbf{B}$ . If  $\mathbf{B}$  is a symmetric PSD matrix, then the eigenvalues of  $\mathbf{B}$  and the singular values of  $\mathbf{B}$  coincide. Let  $\kappa(\mathbf{B}) \stackrel{\text{def}}{=} \sigma_{\max}(\mathbf{B})/\sigma_{\min}(\mathbf{B})$  denote the condition number of  $\mathbf{B}$ .

For two symmetric matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ , we say  $\mathbf{A} \succcurlyeq \mathbf{B}$  iff  $\mathbf{A} - \mathbf{B}$  is PSD. The condition number of  $\mathbf{A}$  relative to  $\mathbf{B}$ , denoted by  $\kappa(\mathbf{A}, \mathbf{B})$ , to be

$$\min \left\{ \frac{\lambda_{\max}}{\lambda_{\min}} : \lambda_{\min} \cdot \mathbf{B} \preccurlyeq \mathbf{A} \preccurlyeq \lambda_{\max} \cdot \mathbf{B} \right\}.$$

Consider the eigen-decomposition of a PSD rank- $k$  matrix  $\mathbf{A} = \sum_{1 \leq i \leq k} \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$ , the *Moore-Penrose pseudo-inverse* of  $\mathbf{A}$  is  $\mathbf{A}^\dagger \stackrel{\text{def}}{=} \sum_{1 \leq i \leq k} \lambda_i^{-1} \mathbf{u}_i \mathbf{u}_i^\top$ .

Given a matrix  $\mathbf{A}$ , let  $\text{im}(\mathbf{A})$  to denote the image of a matrix  $\mathbf{A}$ . Let  $\Pi_{\mathbf{A}} \stackrel{\text{def}}{=} \mathbf{A}(\mathbf{A}\mathbf{A}^\top)^\dagger \mathbf{A}^\top$ ,



i.e. the orthogonal projection onto  $\text{im}(\mathbf{A})$ . Note that  $\mathbf{\Pi}_{\mathbf{A}} = \mathbf{\Pi}_{\mathbf{A}}^\top$  and  $\mathbf{\Pi}_{\mathbf{A}} = \mathbf{\Pi}_{\mathbf{A}}^2$ .

**Norms.** Given a vector  $\mathbf{x} \in \mathbb{R}^n$ , the Euclidean norm (or  $\ell_2$  norm) is defined as  $\|\mathbf{x}\|_2 \stackrel{\text{def}}{=} \sqrt{\sum_{1 \leq i \leq n} \mathbf{x}_i^2}$ . Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a symmetric PSD matrix, we define the matrix norm w.r.t.  $\mathbf{A}$  of a vector  $\mathbf{x}$  as  $\|\mathbf{x}\|_{\mathbf{A}} \stackrel{\text{def}}{=} \sqrt{\mathbf{x}^\top \mathbf{A} \mathbf{x}}$ . For matrix  $\mathbf{A}$ , we define

$$\begin{aligned} \|\mathbf{A}\|_{\max} &\stackrel{\text{def}}{=} \max_{i,j} |\mathbf{A}_{ij}|, \\ \|\mathbf{A}\|_1 &\stackrel{\text{def}}{=} \max_j \sum_i |\mathbf{A}_{ij}|, \\ \|\mathbf{A}\|_\infty &\stackrel{\text{def}}{=} \max_i \sum_j |\mathbf{A}_{ij}|, \\ \|\mathbf{A}\|_2 &\stackrel{\text{def}}{=} \sqrt{\lambda_{\max}(\mathbf{A} \mathbf{A}^\top)}. \end{aligned}$$

We let  $\min_+(\mathbf{A}) = \min_{i,j \text{ s.t. } \mathbf{A}_{ij} \neq 0} |\mathbf{A}_{ij}|$ .

**Schur complement.** Schur complement is an intermediate matrix of partial Gaussian elimination.

**Definition 2.1.1** (Schur complement). *Let  $S, T$  be a partition of the indices of a square matrix  $\mathbf{A}$  so that*

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{SS} & \mathbf{A}_{ST} \\ \mathbf{A}_{ST}^\top & \mathbf{A}_{TT} \end{pmatrix}, \quad (2.1)$$

where  $\mathbf{A}_{SS}, \mathbf{A}_{ST}, \mathbf{A}_{TT}$  are block matrices, the Schur complement of  $\mathbf{A}$  onto  $T$  is

$$\text{SC}[\mathbf{A}]_T \stackrel{\text{def}}{=} \mathbf{A}_{TT} - \mathbf{A}_{ST}^\top \mathbf{A}_{SS}^{-1} \mathbf{A}_{ST}.$$

The following are useful facts of Schur complements, which we will prove in Appendix A.1.

**Fact 2.1.2.** Consider a matrix  $\mathbf{A}$  defined as in Equation (2.1), and any fixed vector  $\mathbf{x}$  of dimension  $|S|$ . Then,

$$\min_{\mathbf{y} \in \mathbb{R}^{|T|}} \begin{pmatrix} \mathbf{x}^\top & \mathbf{y}^\top \end{pmatrix} \mathbf{A} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \mathbf{x}^\top \text{SC}[\mathbf{A}]_T \mathbf{x}.$$

**Fact 2.1.3.** Let  $\mathbf{A}$  be a symmetric PSD matrix and  $\text{SC}[\mathbf{A}]_T$  be its Schur complement.

1.  $\text{SC}[\mathbf{A}]_T$  is a symmetric PSD matrix.
2.  $\lambda_{\max}(\text{SC}[\mathbf{A}]_T) \leq \lambda_{\max}(\mathbf{A})$ .

## 2.2 Approximately Solving A Linear System

In this section we formally define the notions of approximate solutions to linear systems.

**Definition 2.2.1** (Linear System Approximation Problem (LSA)). Given linear system  $(\mathbf{A}, \mathbf{c})$ , where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , and  $\mathbf{c} \in \mathbb{R}^m$ , and given a scalar  $0 \leq \epsilon \leq 1$ , we refer to the LSA problem for the triple  $(\mathbf{A}, \mathbf{c}, \epsilon)$  as the problem of finding  $\mathbf{x} \in \mathbb{R}^n$  s.t.

$$\|\mathbf{A}\mathbf{x} - \mathbf{\Pi}_A \mathbf{c}\|_2 \leq \epsilon \|\mathbf{\Pi}_A \mathbf{c}\|_2,$$

and we say that such an  $\mathbf{x}$  is a solution to the LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$ .

This definition of a LSA instance and solution has several advantages: when  $\text{im}(\mathbf{A}) = \mathbb{R}^m$ , we get  $\mathbf{\Pi}_A = \mathbf{I}$ , and it reduces to the natural condition  $\|\mathbf{A}\mathbf{x} - \mathbf{c}\|_2 \leq \epsilon \|\mathbf{c}\|_2$ , which because  $\text{im}(\mathbf{A}) = \mathbb{R}^m$ , can be satisfied for any  $\epsilon$ , and for  $\epsilon = 0$  tells us that  $\mathbf{A}\mathbf{x} = \mathbf{c}$ .

When  $\text{im}(\mathbf{A})$  does not include all of  $\mathbb{R}^m$ , the vector  $\mathbf{\Pi}_A \mathbf{c}$  is exactly the projection of  $\mathbf{c}$  onto  $\text{im}(\mathbf{A})$ , and so a solution can still be obtained for any  $\epsilon$ . Further, as  $(\mathbf{I} - \mathbf{\Pi}_A)\mathbf{c}$  is

orthogonal to  $\Pi_A \mathbf{c}$  and  $\mathbf{Ax}$ , it follows that

$$\|\mathbf{Ax} - \mathbf{c}\|_2^2 = \|(\mathbf{I} - \Pi_A)\mathbf{c}\|_2^2 + \|\mathbf{Ax} - \Pi_A \mathbf{c}\|_2^2.$$

Thus, when  $\mathbf{x}$  is a solution to the LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$ , then  $\mathbf{x}$  also gives an  $\epsilon^2 \|\Pi_A \mathbf{c}\|_2^2$  additive approximation to

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{c}\|_2^2 = \|(\mathbf{I} - \Pi_A)\mathbf{c}\|_2^2. \quad (2.2)$$

Similarly, an  $\mathbf{x}$  which gives an additive  $\epsilon^2 \|\Pi_A \mathbf{c}\|_2^2$  approximation to Problem (2.2) is always a solution to the LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$ . These observations prove the following (well-known) fact:

**Fact 2.2.2.** *Let  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^m} \|\mathbf{Ax} - \mathbf{c}\|_2^2$ , then for every  $\mathbf{x}$ ,*

$$\|\mathbf{Ax} - \mathbf{c}\|_2^2 \leq \|\mathbf{Ax}^* - \mathbf{c}\|_2^2 + \epsilon^2 \|\Pi_A \mathbf{c}\|_2^2$$

*if and only if  $\mathbf{x}$  is a solution to the LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$ .*

When the linear system  $\mathbf{Ax} = \mathbf{c}$  does not have a solution, a natural notion of solution is any minimizer of Problem (2.2). A simple calculation shows that this is equivalent to requiring that  $\mathbf{x}$  is a solution to the linear system  $\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{c}$ , which always has a solution even when  $\mathbf{Ax} = \mathbf{c}$  does not. The system  $\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{c}$  is referred to as the *normal equation* associated with  $\mathbf{Ax} = \mathbf{c}$  (see [62]).

**Fact 2.2.3.**  *$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{c}\|_2^2$ , if and only if  $\mathbf{A}^\top \mathbf{Ax}^* = \mathbf{A}^\top \mathbf{c}$ , and this linear system always has a solution.*

This leads to a natural question: Suppose we want to approximately solve the linear system  $\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{c}$ . Can we choose our notion of approximation to be equivalent to that of a solution to the LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$ ?

A second natural question is whether we can choose a notion of distance between a proposed solution  $\mathbf{x}$  and an optimal solution  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{A}\mathbf{x} - \mathbf{c}\|_2^2$  s.t. this distance being small is equivalent to  $\mathbf{x}$  being a solution to the LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$ ? The answer to both questions is yes, as demonstrated by the following facts:

**Fact 2.2.4.** *Suppose  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{A}\mathbf{x} - \mathbf{c}\|_2^2$  then*

1.  $\|\mathbf{A}^\top \mathbf{A}\mathbf{x} - \mathbf{A}^\top \mathbf{c}\|_{(\mathbf{A}^\top \mathbf{A})^\dagger} = \|\mathbf{A}\mathbf{x} - \Pi_{\mathbf{A}} \mathbf{c}\|_2 = \|\mathbf{x} - \mathbf{x}^*\|_{\mathbf{A}^\top \mathbf{A}}.$
2. *The following statements are each equivalent to  $\mathbf{x}$  being a solution to the LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$ :*
  - (a)  $\|\mathbf{A}^\top \mathbf{A}\mathbf{x} - \mathbf{A}^\top \mathbf{c}\|_{(\mathbf{A}^\top \mathbf{A})^\dagger} \leq \epsilon \|\mathbf{A}^\top \mathbf{c}\|_{(\mathbf{A}^\top \mathbf{A})^\dagger}$  *if and only if  $\mathbf{x}$  is a solution to the LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$ .*
  - (b)  $\|\mathbf{x} - \mathbf{x}^*\|_{\mathbf{A}^\top \mathbf{A}} \leq \epsilon \|\mathbf{x}^*\|_{\mathbf{A}^\top \mathbf{A}}$  *if and only if  $\mathbf{x}$  is a solution to the LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$ .*

For completeness, we prove Fact 2.2.4 in Appendix A.2. Fact 2.2.4 explains connection between our Definition 2.2.1, and the usual convention for measuring error in the Laplacian solver literature [9]. In this setting, we consider a Laplacian matrix  $\mathbf{L}$ , which can be written as  $\mathbf{L} = \mathbf{A}^\top \mathbf{A} \in \mathbb{R}^{n \times n}$ , and a vector  $\mathbf{b}$  s.t.  $\Pi_{\mathbf{A}^\top \mathbf{A}} \mathbf{b} = \mathbf{b}$ . This condition on  $\mathbf{b}$  is easy to verify in the case of Laplacians, since for the Laplacian of a connected graph,  $\Pi_{\mathbf{A}^\top \mathbf{A}} = \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top$ . Additionally, it is also equivalent to the condition that there exists  $\mathbf{c}$  s.t.  $\mathbf{b} = \mathbf{A}^\top \mathbf{c}$ . For Laplacians it is possible to compute both  $\mathbf{A}$  and a vector  $\mathbf{c}$  s.t.  $\mathbf{b} = \mathbf{A}^\top \mathbf{c}$  in time linear in  $\text{nnz}(\mathbf{L})$ . For Laplacian solvers, the approximation error of an approximate solution  $\mathbf{x}$  is measured by the  $\epsilon$  s.t.  $\|\mathbf{A}^\top \mathbf{A}\mathbf{x} - \mathbf{b}\|_{(\mathbf{A}^\top \mathbf{A})^\dagger} \leq \epsilon \|\mathbf{b}\|_{(\mathbf{A}^\top \mathbf{A})^\dagger}$ . By Fact 2.2.4, we see that this is exactly equivalent to  $\mathbf{x}$  being a solution to the LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$ .

The following Lemma states that, if  $\Pi_{\mathbf{A}} \mathbf{c} \neq \mathbf{0}$ , then its norm is separated from 0. This lemma is important in our proofs of hardness results for structured linear systems and packing LPs. We will prove it in Appendix A.2.

**Lemma 2.2.5.** *Let  $\mathbf{A}$  and  $\mathbf{c}$  a matrix and a vector such that  $\|\mathbf{A}\mathbf{c}\|_2^2$  is integral. If  $\mathbf{A}^\top \mathbf{c} \neq \mathbf{0}$ , then*

$$\|\Pi_{\mathbf{A}} \mathbf{c}\|_2^2 \geq \frac{1}{\sigma_{\max}^2(\mathbf{A})} = \frac{1}{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})}.$$

**Remark.** We assume that in Problem (2.2)  $\mathbf{A}$  in general has dimension  $\tilde{O}(n) \times n$ , otherwise we can use subspace embedding [63] to compute an  $\tilde{\mathbf{A}} \in \mathbb{R}^{\tilde{O}(n) \times n}$  which approximately preserves the  $\ell_2$  norm of  $\mathbf{A}\mathbf{x}$  for any  $\mathbf{x}$ . In addition, computing  $\tilde{\mathbf{A}}$  takes  $\text{nnz}(\mathbf{A})$  time, and it preserves the sparsity of  $\mathbf{A}$ .

### 2.3 Graph-Structured Block Matrices

In this subsection, we formally define Graph-Structured Block Matrices, and some matrices belonging to this class: 2-commodity Laplacian matrices, truss stiffness matrices, and 2 total variation matrices.

Suppose we have a collection of  $n$  disjoint sets of variables  $X_1, \dots, X_n$ , with each set having the same size,  $|X_i| = d$ . Let  $\mathbf{x}^i$  denote the vector of variables in  $X_i$ , and consider an equation of the form  $\mathbf{M}\mathbf{x}^i - \mathbf{N}\mathbf{x}^j = 0$ , where  $\mathbf{M}$  and  $\mathbf{N}$  are both  $r \times d$  matrices. Now we form a linear system  $\mathbf{B}\mathbf{x} = \mathbf{0}$  by stacking  $m$  equations of the form given above as the rows of the system. Note that, very importantly, we allow a different choice of  $\mathbf{M}$  and  $\mathbf{N}$  for every pair of  $i$  and  $j$ . This matrix  $\mathbf{B} \in \mathbb{R}^{mr \times nd}$  we refer to as an *Incidence-Structured Block Matrix* (ISBM), while we refer to  $\mathbf{B}^\top \mathbf{B}$  as a *Graph-Structured Block Matrix* (GSBM). Note that  $\mathbf{B}$  is not usually PSD, but  $\mathbf{B}^\top \mathbf{B}$  is. The number of non-zeros in  $\mathbf{B}^\top \mathbf{B}$  is  $O(md^2)$ . GSBMs come up in many applications, where we typically want to solve a linear system in the normal equations of  $\mathbf{B}$ .

Laplacian matrices are GSBMs where  $d = 1$  and  $\mathbf{M} = \mathbf{N} = w$ , where  $w$  is a real number, and we allow different  $w$  for each pair of  $i$  and  $j$ . The corresponding ISBM for Laplacians is called an edge-vertex incidence matrix. Connection Laplacians are GSBMs where  $d = O(1)$  and  $\mathbf{M} = \mathbf{N}^\top = w\mathbf{Q}$ , for some rotation matrix  $\mathbf{Q}$  and a real num-

ber  $w$ . Again, we allow a different rotation matrix and scaling for every edge. For both Laplacians and Connection Laplacians, there exist linear system solvers that run in time  $O(m \text{ polylog}(n, \epsilon^{-1}))$  and produce  $\epsilon$  approximate solutions to the corresponding normal equations.

We now introduce several classes of ISBMs and their associated GSBMs.

**Definition 2.3.1** (2-commodity incidence matrix). *A 2-commodity incidence matrix is an ISBM where  $d = 2$  and  $r = 1$ , and  $\mathbf{M} = \mathbf{N}$ , and we allow three types of  $\mathbf{M}$ :  $\mathbf{M} = w \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$ ,  $\mathbf{M} = w \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$  and  $\mathbf{M} = w \begin{pmatrix} 1 & -1 \\ 0 & 0 \end{pmatrix}$ , where in each case  $w$  is a real number which may depend on the pair  $i$  and  $j$ . We denote the set of all 2-commodity incidence matrices by  $\mathcal{MC}_2$ . The corresponding GSBM is called a 2-commodity Laplacian. The ISBM definition is equivalent to requiring the GSBM to have the form*

$$\mathbf{L}^1 \otimes \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + \mathbf{L}^2 \otimes \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} + \mathbf{L}^{1+2} \otimes \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$$

where  $\otimes$  is the tensor product and  $\mathbf{L}^1$ ,  $\mathbf{L}^2$ , and  $\mathbf{L}^{1+2}$  are all Laplacian matrices.

We adopt a convention that the first variable in a set  $X_i$  is labelled  $\mathbf{u}_i$  and the second is labelled  $\mathbf{v}_i$ . Using this convention, given a 2-commodity incidence matrix  $\mathbf{B}$ , the equation  $\mathbf{B}\mathbf{x} = \mathbf{0}$  must consist of scalings of the following three types of equations:  $\mathbf{u}_i - \mathbf{u}_j = 0$ ,  $\mathbf{v}_i - \mathbf{v}_j = 0$ , and  $\mathbf{u}_i - \mathbf{v}_i - (\mathbf{u}_j - \mathbf{v}_j) = 0$ .

**Definition 2.3.2** (Strict 2-commodity incidence matrix). *A strict 2-commodity incidence matrix is a 2-commodity incidence matrix where the corresponding 2-commodity Laplacian has the property that  $\mathbf{L}^1$ ,  $\mathbf{L}^2$ , and  $\mathbf{L}^{1+2}$  all have the same non-zero pattern. We denote the set of all strict 2-commodity incidence matrices by  $\mathcal{MC}_2^{>0}$ . We denote the set of all strict 2-commodity incidence matrices with integer entries by  $\mathcal{MC}_{2,\mathbb{Z}}^{>0}$ .*

Linear systems in  $\mathcal{MC}_2^{>0}$  are exactly the systems that one has to solve to when solving 2-commodity problems using Interior Point Methods (IPMs). For readers unfamiliar with

2-commodity problems or IPMs, we provide a brief explanation of why this is the case in Section 3.8. The  $\mathcal{MC}_2^{>0}$  is more restrictive than  $\mathcal{MC}_2$ , and  $\mathcal{MC}_{2,\mathbb{Z}}^{>0}$  in turn is even more restrictive. One could hope that fast linear system solvers exist for  $\mathcal{MC}_2^{>0}$  or  $\mathcal{MC}_{2,\mathbb{Z}}^{>0}$ , even if they do not exist for  $\mathcal{MC}_2$ . However, our reductions show that even getting a fast approximate solver for  $\mathcal{MC}_{2,\mathbb{Z}}^{>0}$  with polynomially bounded entries and condition number will lead to a fast solver for all matrices with polynomially bounded entries and condition number.

The next class we consider is Truss Stiffness Matrices. They have been studied extensively in the numerical linear algebra community [43, 64].

**Definition 2.3.3** (2D Truss Incidence Matrices). *Let  $G = (V, E)$  be a graph whose vertices are  $n$  points in 2-dimension:  $s^1, \dots, s^n \in \mathbb{R}^2$ . Consider  $X_1, \dots, X_n$  where  $d = 2$ . A 2D Truss Incidence Matrix is an ISBM where  $d = 2$  and  $r = 1$ , and for each  $i$  and  $j$ ,  $M = w(s^i - s^j)^\top$ , and  $w$  is a real number that may depend on the pair  $i$  and  $j$ , but  $s^i$  depends only on  $i$  and vice versa for  $s^j$ . We denote the class of all 2D Truss Incidence Matrices by  $\mathcal{T}_2$ .*

In Section 2.4, we extend the above definition for 2D trusses into 3D.

Another important class of matrices is Total Variation Matrices (TV matrices). TV matrices come from Interior Point Methods for solving total variation minimization problem in image, see for example [65] and [32]. Not all TV matrices are GSBMs, but many GSBMs can be expressed as TV matrices.

**Definition 2.3.4** (TV matrix and 2-TV Incidence Matrices). *Let  $E_1 \cup \dots \cup E_s$  be a partition of the edge set of a graph. For each  $1 \leq i \leq s$ , let  $B^i$  be the edge-vertex incidence matrix of  $E_i$ ,  $W^i$  be a diagonal matrix of edge weights, and  $r^i$  be a vector satisfying  $W^i \succcurlyeq r^i(r^i)^\top$ . The total variation matrix (TV matrix) is defined as*

$$M = \sum_{1 \leq i \leq s} (B^i)^\top (W^i - r^i(r^i)^\top) B^i.$$

A 2-TV Incidence Matrix is defined as any ISBM whose corresponding GSBM is a TV

matrix with  $\mathbf{W}^i \in \mathbb{R}^{2 \times 2}$  and  $\mathbf{r}^i \in \mathbb{R}^2$ . We denote the class of all 2-TV incidence matrices by  $\mathcal{V}_2$ .

We remark that, given an Incidence-Structured Block Matrix  $\mathbf{B} \in \mathbb{R}^{mr \times nd}$  with  $r, d$  being constant (this is usually the case), we can compute its Graph-Structured Block Matrix  $\mathbf{B}^\top \mathbf{B}$  in  $\text{nnz}(\mathbf{B})$  time. Conversely, given a GSBM  $\mathbf{M}$ , we can decompose it to  $\mathbf{M} = \mathbf{B}^\top \mathbf{B}$  for an ISBM  $\mathbf{B}$  in  $\text{nnz}(\mathbf{M})$  time. Thus, transforming between ISBMs and GSBMs is never a bottleneck of running time.

## 2.4 Geometric Structures for 3-D Trusses

In this section, we give more details about the class of 3-D trusses which we will discuss in Section 5.

**Tetrahedral meshes.** For a subset  $S \subset \mathbb{R}^3$ , we define the *diameter* of  $S$  to be the maximum Euclidean distance between any pair of points in  $S$ . We define the *aspect ratio* of  $S$  to be the ratio between the radius of the smallest ball containing  $S$  and the radius of the largest ball inscribed in  $S$ . In mesh generation, aspect ratio is a common criterion for individual elements.

A tetrahedron is the convex hull of four non-coplanar points in  $\mathbb{R}^3$ . We will specify tetrahedrons in terms of sets of four such points. We refer to a set of tetrahedrons as a tetrahedral mesh.

We say a tetrahedral mesh is a *simplicial complex* if the intersection of every two tetrahedrons is either empty or a face of both two tetrahedrons. A simplicial complex is *convex* if the union of the images of its simplices is convex, as defined in [66].

The following definition characterizes rigidity and stiffness of a tetrahedral mesh, which is an adaption of Definition 2.3 in [44].

**Definition 2.4.1.** *The rigidity graph of a tetrahedral mesh is a graph whose vertices correspond to tetrahedrons and whose edges connect any two tetrahedrons sharing a triangle*



face. A tetrahedral mesh is stiffly-connected iff (1) its rigidity graph is connected, and (2) for any vertex in the mesh, the rigidity subgraph induced on the tetrahedrons containing this vertex is connected.

We define a *bounding box* of a convex 3D shape to be a 3D box such that: (1) the box contains all points of this shape, and (2) the volume of the box is same as the volume of this shape up to a constant factor. [67] gives a linear time algorithm which computes a bounding box of a convex shape in 3 dimensions.

**Lemma 2.4.2.** (Lemma 3.6 of [67]) *Given a 3D convex shape, one can compute in linear time a bounding box of this shape.<sup>1</sup> Moreover, the aspect ratio of the bounding box is same as the aspect ratio of the given shape up to a constant factor.*

**3D trusses and stiffness matrices.** We use the following notations for 3D trusses.

**Definition 2.4.3** (3-dimensional truss). 3-dimensional truss  $\mathcal{T} = \langle V, \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$  is given by

- A set of  $n$  vertices  $V$  embedded at distinct points  $\mathbf{p}_1, \dots, \mathbf{p}_n \in \mathbb{R}^3$ .
- A mesh (i.e. set) of tetrahedrons  $T = \{t_1, t_2, \dots\}$ , each specified in terms of four vertices, i.e. we identify tetrahedron  $t_i$  with both four vertices  $\{a_i, b_i, c_i, d_i\} \subseteq V$  and the convex hull of  $\mathbf{p}_{a_i}, \mathbf{p}_{b_i}, \mathbf{p}_{c_i}, \mathbf{p}_{d_i}$ .
- A set of edges  $E$  which is exactly the set of pairs of vertices that appear in some tetrahedron together. Each edge  $e = (i, j) \in E$  represents a straight idealized bar between vertex points  $\mathbf{p}_i$  and  $\mathbf{p}_j$ .
- A function  $\gamma : E \rightarrow \mathbb{R}_+$ , which assigns a stiffness coefficient  $\gamma(e)$  to each edge  $e$ . The stiffness coefficient represents the stiffness of the idealized bar corresponding to edge  $e$ .

---

<sup>1</sup>The lemma statement in [67] gives a bound related to the minimum-volume bounding box  $B^*$  of the input shape, but their proof uses the volume of the shape as a lower bound of the volume of  $B^*$ .

The following definition is a restatement and generalization of Definition 2.3.3. Definition 2.4.4 is more useful in our algorithm design in Section 5.

**Definition 2.4.4** (3-D Truss stiffness matrix). *Let  $\mathcal{T} = \langle V, \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$  be a 3-dimensional truss. For each edge  $e = (i, j) \in E$ , we define an edge vector  $\mathbf{b}^{(e)} \in \mathbb{R}^{3n}$  with 6 nonzero entries:*

$$\mathbf{b}_{3i-2:3i}^{(e)} = -\mathbf{b}_{3j-2:3j}^{(e)} = \frac{\mathbf{p}_i - \mathbf{p}_j}{\|\mathbf{p}_i - \mathbf{p}_j\|_2}.$$

*The stiffness matrix of the truss  $\mathcal{T}$  is defined as*

$$\mathbf{A}_{\mathcal{T}} \stackrel{\text{def}}{=} \sum_{e=(i,j) \in E} \frac{\gamma(e)}{\|\mathbf{p}_i - \mathbf{p}_j\|_2} \mathbf{b}^{(e)} \mathbf{b}^{(e)\top}.$$

In addition, we use some terms to simplify our descriptions.

**Definition 2.4.5.** *A 3D truss is said to be simple iff it is a simplicial complex and every tetrahedron has bounded aspect ratio. We say a 3D truss is edge-simple iff its tetrahedral mesh is simple and every tetrahedron has bounded edge lengths and stiffness coefficients (i.e. both are bounded above and below by constants). We say a 3D truss is convex edge-simple iff it is edge-simple, and its tetrahedral mesh is convex.*

## 2.5 Techniques for Solving Linear Systems

Our algorithm combines two of the most important tools for solving linear systems: Nested dissection and preconditioning. Below, we give a brief introduction to some of the central results on these techniques.

Classic results due to Lipton, Rose, and Tarjan [48], and Miller and Thurston [49] combine to show that linear systems arising from simple tetrahedral meshes (see Definition 2.4.5) can be solved in  $O(n^2)$  time. These results concern linear equations in an  $n \times n$  matrix  $\mathbf{A}$  where the indices  $\{1, \dots, n\}$  can be embedded as points  $\{\mathbf{p}_1, \dots, \mathbf{p}_n\}$  that form

the vertices of an explicitly given, simple tetrahedral mesh, and  $A_{ij}$  is non-zero only if the vertices  $i$  and  $j$  share an edge in the tetrahedral mesh.

**Theorem 2.5.1** (Nested dissection [49]). *Let  $A \in \mathbb{R}^{n \times n}$  be a symmetric matrix defined on a simple tetrahedral mesh. A Cholesky factorization  $A = PLL^\top P^\top$  can be computed in time  $O(n^2)$ , in which  $P$  is a permutation matrix and  $L$  is a lower triangular matrix with  $O(n^{4/3})$  nonzero entries. As a result, a linear system in  $A$  can be solved in time  $O(n^2)$  by Gaussian elimination.*

Theorem 2.5.1 can be extended to a block matrix  $A \in \mathbb{R}^{cn \times cn}$  where  $c$  is a constant positive integer. Each vertex of the underlying graph corresponds to  $c$  indices of  $A$ . In addition, the block corresponding to the column indices for vertex  $i$  and the row indices for  $j$  should be non-zero only if the vertices  $i$  and  $j$  share an edge in the tetrahedral mesh, or if  $i = j$ , i.e. when the block is on the diagonal.

The Nested Dissection algorithm relies on invoking *separators* recursively. A separator is a set of indices  $S$  such that the remaining indices  $[n] \setminus S$  can be partitioned into two sets  $B$  and  $C$  such that every entry with  $i \in B$  and  $j \in C$  has  $A_{ij} = 0$ . Furthermore, we guarantee that the partition is roughly balanced, for example, each of  $B$  and  $C$  contains no more than  $\frac{3}{4} \cdot n$  indices. Nested Dissection recursively repeats the partitioning process on the union of each subset and the separator itself, that is,  $B \cup S$  and  $C \cup S$ . Given such a recursive partition scheme, we reorder the indices of the matrix so that the indices in the separator  $S$  are eliminated *last*, and we then order the indices in  $B$  and  $C$  recursively in a similar way. We perform Gaussian elimination on the matrix according to this ordering, which only introduces a small fill-in size and few multiplication counts. This approach also works for eliminating a subset of the variables, resulting in a Schur complement on the rest.

Both the running time and representation cost of nested dissection algorithms are bottlenecked by the costs of the top-level separators. In Algorithm 12 TRUSSSOLVER, we will utilize improved running time bounds for nested dissection when better separators exist. The following lemma characterizes the performance of Nested Dissection given better

top-level separators.

**Lemma 2.5.2.** *Suppose we have a recursive separator decomposition of a simplicial complex with  $n$  bounded aspect ratio tetrahedrons such that:*

1. *the number of leaves, and hence total number of recursive calls, is at most  $n^\alpha$ .*
2. *each leaf (bottom layer partition) has at most  $n^\beta$  tetrahedrons.*
3. *each top separator has size at most  $n^\gamma$ .*

*Then we can find an exact Cholesky factorization of the associated stiffness matrix in time  $O(n^{\alpha+2\beta} + n^{\alpha+3\gamma})$ , and the total resulting fill-in is  $O(n^{\alpha+\frac{4}{3}\beta} + n^{\alpha+2\gamma})$ .*

We give a proof of the above lemma in Section 5.5.1, which is an adaption of the analysis in [48] We remark that the algorithmic realization of this can be viewed as utilizing the nested dissection algorithm 2.5.1 to complete this structure into a full separator tree.

Last but not least, we state the following theorem for preconditioned conjugate gradient, which will be used in bounding the running time of our algorithm.

**Theorem 2.5.3** (Preconditioned conjugate gradient [68]). *Let  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$  be two symmetric positive semidefinite matrices and let  $\mathbf{b} \in \mathbb{R}^n$ . Each iteration of the preconditioned conjugate gradient multiplies one vector by  $\mathbf{A}$ , solves one linear system in  $\mathbf{B}$ , and performs a constant number of vector additions. For any  $\epsilon > 0$ , the algorithm outputs an  $\mathbf{x}$  satisfying  $\|\mathbf{Ax} - \mathbf{b}\|_2 \leq \epsilon \|\mathbf{b}\|_2$  in  $O(\sqrt{\kappa(\mathbf{A}, \mathbf{B})} \log(1/\epsilon))$  such iterations.*

We remark that while there are settings where the convergence of preconditioned conjugate gradient is numerically unstable, the eigenvalue-based bound that we utilize here is stable once the solves involving  $\mathbf{B}$  have polynomially small errors.

## CHAPTER 3

### HARDNESS RESULTS OF STRUCTURED LINEAR SYSTEMS

In this chapter, we prove Theorem 1.1.1.

#### 3.1 Measuring the difficulty of solving a linear system

Running times for iterative linear system solvers generally depend on the number of non-zeros in the input matrix, the condition number of the input matrix, the accuracy, and the bit complexity.

In this section, we formally define several measures of complexity of the linear systems we use. This is crucial, because we want to make sure that our reductions do not rely on mapping into extremely ill-conditioned matrices, and so we use these measures to show that this is in fact not the case.

**Definition 3.1.1.** *The sparse parameter complexity of an LSA instance  $(\mathbf{A}, \mathbf{c}, \epsilon)$  where  $\mathbf{A} \in \mathbb{Z}^{m \times n}$  and  $\text{nnz}(\mathbf{A}) \geq \max(m, n)$ , and  $\epsilon > 0$ , is*

$$\mathcal{S}(\mathbf{A}, \mathbf{c}, \epsilon) \stackrel{\text{def}}{=} \left( \text{nnz}(\mathbf{A}), \max \left( \|\mathbf{A}\|_{\max}, \|\mathbf{c}\|_{\max}, \frac{1}{\min_+(\mathbf{A})}, \frac{1}{\min_+(\mathbf{c})} \right), \kappa(\mathbf{A}), \epsilon^{-1} \right).$$

Note in the definition above that when  $\mathbf{A} \neq \mathbf{0}$  and  $\mathbf{c} \neq \mathbf{0}$  have only integer entries, we trivially have  $\min_+(\mathbf{A}) \geq 1$  and  $\min_+(\mathbf{c}) \geq 1$ . However, including  $\frac{1}{\min_+(\mathbf{A})}$ , and  $\frac{1}{\min_+(\mathbf{c})}$  in the definition stated above is useful when working with intermediate matrices whose entries are not integer valued.

The following Claim bounds some frequently used quantities by parameters defined in Definition 3.1.1. We prove it in Appendix A.2.

**Claim 3.1.2.** *Given an LSA  $(\mathbf{A}, \mathbf{c}, \epsilon)$  satisfying:  $\mathbf{A}$  and  $\mathbf{c}$  are integral,  $\mathbf{A}$  has  $s$  non-zeros and condition number  $K$ , and every entry of  $\mathbf{A}$  and  $\mathbf{c}$  has absolute value at most  $U$ . Then*

1.  $\|\mathbf{A}\|_\infty \leq sU$ .
2.  $\|\mathbf{c}\|_2 \leq \sqrt{s}U$ .
3.  $\frac{1}{sU^2} \leq \lambda_{\max}(\mathbf{A}^\top \mathbf{A}) \leq sU^2$  and  $\lambda_{\min}(\mathbf{A}^\top \mathbf{A}) \geq \frac{1}{sK^2U^2}$ .

We use the term *matrix class* to refer to an infinite set of matrices  $\mathcal{M}$ . In this section, we formally define a notion of efficient reduction between linear systems in different classes of matrices.

**Definition 3.1.3** (Efficient  $f$ -reducibility). *Suppose we have two matrix classes  $\mathcal{M}_1$  and  $\mathcal{M}_2$ , and there exist two algorithms  $\mathcal{A}_{1 \rightarrow 2}$  and  $\mathcal{A}_{1 \leftarrow 2}$  s.t. given an LSA instance  $(\mathbf{M}^1, \mathbf{c}^1, \epsilon)$ , where  $\mathbf{M}^1 \in \mathcal{M}_1$ , the call  $\mathcal{A}_{1 \rightarrow 2}(\mathbf{M}^1, \mathbf{c}^1, \epsilon_1)$  returns an LSA instance  $(\mathbf{M}^2, \mathbf{c}^2, \epsilon_2)$  s.t. if  $\mathbf{x}^2$  is a solution to LSA instance  $(\mathbf{M}^2, \mathbf{c}^2, \epsilon_2)$  then  $\mathbf{x}^1 = \mathcal{A}_{1 \leftarrow 2}(\mathbf{M}^1, \mathbf{M}^2, \mathbf{x}^2)$  is a solution to LSA instance  $(\mathbf{M}^1, \mathbf{c}^1, \epsilon_1)$ .*

*Consider a function of  $f : \mathbb{R}_+^4 \rightarrow \mathbb{R}_+^4$  s.t. every output coordinate is an increasing function of every input coordinate. Suppose that we always have*

$$\mathcal{S}(\mathbf{M}^2, \mathbf{c}^2, \epsilon_2) \leq f(\mathcal{S}(\mathbf{M}^1, \mathbf{c}^1, \epsilon_1)),$$

*and the running times of  $\mathcal{A}_{1 \rightarrow 2}(\mathbf{M}^1, \mathbf{c}^1, \epsilon_1)$  and  $\mathcal{A}_{1 \leftarrow 2}(\mathbf{M}^1, \mathbf{M}^2, \mathbf{x}^2)$  are both bounded by  $O(\text{nnz}(\mathbf{M}^1))$ .*

*Then we say that  $\mathcal{M}_1$  is efficiently  $f$ -reducible to  $\mathcal{M}_2$ , which we also write as*

$$\mathcal{M}_1 \leq_f \mathcal{M}_2.$$

**Lemma 3.1.4.** *Suppose  $\mathcal{M}_1 \leq_f \mathcal{M}_2$  and  $\mathcal{M}_2 \leq_g \mathcal{M}_3$ . Then  $\mathcal{M}_1 \leq_{g \circ f} \mathcal{M}_3$ .*

*Proof.* The proof is simply by the trivial composition of the two reductions. □

**Definition 3.1.5.** We let  $\mathcal{G}$  denote the class of all matrices with integer valued entries s.t. there is at least one non-zero entry in every row and column<sup>1</sup>.

### 3.2 Main Results and Reduction Outline

In this section, we use the notions of sparse parameter complexity and matrix class reductions defined in Section 3.1 to prove our main technical result, Theorem 3.2.1. Then, as a corollary, we get Theorem 3.2.2.

**Theorem 3.2.1.** Let  $f(s, U, K, \epsilon) = (O(s \log(sU)), \text{poly}(UK\epsilon^{-1}s), \text{poly}(UK\epsilon^{-1}s), \text{poly}(UK\epsilon^{-1}s))$ , then

1.  $\mathcal{G} \leq_f \mathcal{MC}_{2,\mathbb{Z}}^{>0}$ .
2.  $\mathcal{G} \leq_f \mathcal{T}_2$ .
3.  $\mathcal{G} \leq_f \mathcal{V}_2$ .

**Theorem 3.2.2.** Suppose we have an algorithm which solves every Linear System Approximation Problem  $(\mathbf{A}, \mathbf{c}, \epsilon)$  with sparse parameter complexity  $\mathcal{S}(\mathbf{A}, \mathbf{c}, \epsilon) \leq (s, U, K, \epsilon^{-1}) O(s^a \text{polylog}(s, U, K, \epsilon^{-1}))$  for some  $a \geq 1$ , whenever  $\mathbf{A} \in \mathcal{R}$  for at least one of  $\mathcal{R} \in \{\mathcal{MC}_{2,\mathbb{Z}}^{>0}, \mathcal{T}_2, \mathcal{V}_2\}$ . I.e. we have a “fast” solver<sup>2</sup> for one of the matrix classes  $\mathcal{MC}_{2,\mathbb{Z}}^{>0}, \mathcal{T}_2$ , or  $\mathcal{V}_2$ . Then every Linear System Approximation Problem  $(\mathbf{A}, \mathbf{c}, \epsilon)$  where  $\mathbf{A} \in \mathcal{G}$  with sparse parameter complexity  $\mathcal{S}(\mathbf{A}, \mathbf{c}, \epsilon) \leq (s, U, K, \epsilon^{-1})$  can be solved in time  $O(s^a \text{polylog}(s, U, K, \epsilon^{-1}))$ .

*Proof.* The theorem is a immediate corollary of Theorem 3.2.1. □

**Definition 3.2.3.** We let  $\mathcal{G}_{z,2}$  denote the class of all matrices with integer valued entries s.t. there is at least one non-zero entry in every row and column, and every row has zero row sum, and for each row, the sum of the positive coefficients is a power of 2.

<sup>1</sup>If there is a row or column with only zeros, then it can always be handled trivially in the context of solving linear systems

<sup>2</sup>The reduction requires only a single linear system solve, and uses the solution in a black-box way. So the reduction also applies if the solver for the class  $\mathcal{R}$  only works with high probability or only has running time guarantees in expectation.

Our reduction follows the steps:

$$\mathcal{G} \rightarrow \mathcal{G}_{z,2} \rightarrow \mathcal{MC}_2 \rightarrow \mathcal{MC}_2^{>0} \rightarrow \mathcal{MC}_{2,\mathbb{Z}}^{>0}.$$

Recall that, from Definition 2.3.2,  $\mathcal{MC}_2^{>0}$  is the class of strict 2-commodity matrices, and  $\mathcal{MC}_{2,\mathbb{Z}}^{>0}$  is the class of strict 2-commodity matrices with integral entries. Specifically, we will prove the following sequence of lemmas.

**Lemma 3.2.4.** *Let  $f(s, U, K, \epsilon) = (O(s), O(\epsilon^{-1}s^{9/2}U^3), O(\epsilon^{-1}s^8U^3K), O(s^{5/2}U^2\epsilon^{-1}))$ , then*

$$\mathcal{G} \leq_f \mathcal{G}_{z,2}.$$

**Lemma 3.2.5.** *Let  $f(s, U, K, \epsilon) = (O(s \log(sU)), O(s^{3/2}U \log^{1/2}(sU)), O(Ks^4U^2 \log^2(sU)), O(sU^2\epsilon^{-1}))$ , then*

$$\mathcal{G}_{z,2} \leq_f \mathcal{MC}_2.$$

**Lemma 3.2.6.** *Let  $f(s, U, K, \epsilon) = (O(s), O(\epsilon^{-1}U^2K), O(\epsilon^{-1}s^2U^2K), O(\epsilon^{-1}))$ , then*

$$\mathcal{MC}_2 \leq_f \mathcal{MC}_2^{>0}.$$

**Lemma 3.2.7.** *Let  $f(s, U, K, \epsilon) = (s, \epsilon^{-1}sU, 2K, O(\epsilon^{-1}))$ , then*

$$\mathcal{MC}_2^{>0} \leq_f \mathcal{MC}_{2,\mathbb{Z}}^{>0}.$$

**Lemma 3.2.8.** *Let  $f(s, U, K, \epsilon)$  be as defined in Lemma 3.2.5 then*

$$\mathcal{G}_{z,2} \leq_f \mathcal{T}_2.$$



**Lemma 3.2.9.** *Let  $f(s, U, K, \epsilon) = (s, U, K, \epsilon^{-1})$ , then*

$$\mathcal{MC}_2 \leq_f \mathcal{V}_2.$$

*Proof of Theorem 3.2.1.* Follows by appropriate composition (Lemma 3.1.4) applied to the the Lemmas above, i.e. 3.2.4, 3.2.5, 3.2.6, 3.2.7, 3.2.8 and 3.2.9.  $\square$

**Outline of the Remaining Section.** We prove Lemma 3.2.4 in Section 3.6, Lemma 3.2.5 in Section 3.3, Lemma 3.2.6 in Section 3.4, Lemma 3.2.7 in Section 3.5, Lemma 3.2.8 in Section 3.7, and Lemma 3.2.9 in Section 3.8.2. In addition, we show how to derive multi-commodity Laplacians and Total Variation matrices by Interior-point method in Section 3.8.

### 3.3 Reducing Zero-Sum Power Two Linear Systems to Two-Commodity Linear Systems

To prove Lemma 3.2.5, we need to provide mapping algorithms  $\mathcal{A}_{\mathcal{G}_{z,2} \rightarrow \mathcal{MC}_2}$  for mapping linear system approximation (LSA) instances over matrices in  $\mathcal{G}_{z,2}$  to LSA problem instances over matrices in  $\mathcal{MC}_2$ , as well as  $\mathcal{A}_{\mathcal{MC}_2 \leftarrow \mathcal{G}_{z,2}}$  for mapping the resulting solutions back. These leads to the following main components:

- Algorithm 1 states the pseudo-code for the algorithm  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$ , which implements the desired mapping of problem instances. Given LSA problem instance  $(\mathbf{A}, \mathbf{c}^A, \epsilon^A)$  where  $\mathbf{A} \in \mathcal{G}_{z,2}$ , the call  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2(\mathbf{A}, \mathbf{c}^A, \epsilon^A)$  returns an LSA problem instance  $(\mathbf{B}, \mathbf{c}^B, \epsilon^B)$  where  $\mathbf{B} \in \mathcal{MC}_2$ . (Strictly speaking,  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$  also has a parameter  $\alpha$  which we will set before using the algorithm.)
- Algorithm 2 provides the pseudo-code for  $\mathcal{MC}_2\text{Gadget}$ , a short subroutine used in  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$  to represent the equation  $a + b = 2c$  using two-commodity

constraints.

- Algorithm 3 provides the (trivial) pseudo-code for  $\text{MAPSOLN } \mathcal{MC}_2 \text{ TO } \mathcal{G}_{z,2}$  used to map solutions to LSA problems over  $\mathcal{MC}_2$  back to solutions over  $\mathcal{G}_{z,2}$  by restricting onto the original variables.

Pseudocode of the key reduction routine that creates the new linear system,  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$ , is shown in Algorithm 1. Note that given two singleton multi-sets each containing a single equation, e.g.  $\{\mathbf{a}_1^\top \mathbf{x} = c_1\}$  and  $\{\mathbf{a}_2^\top \mathbf{x} = c_2\}$  where  $\mathbf{a}_1, \mathbf{a}_2$  are vectors, we define  $\{\mathbf{a}_1^\top \mathbf{x} = c_1\} + \{\mathbf{a}_2^\top \mathbf{x} = c_2\} = \{\mathbf{a}_1^\top \mathbf{x} + \mathbf{a}_2^\top \mathbf{x} = c_1 + c_2\}$  and we define  $\{\mathbf{a}_1^\top \mathbf{x} = c_1\} \cup \{\mathbf{a}_2^\top \mathbf{x} = c_2\} = \{\mathbf{a}_1^\top \mathbf{x} = c_1, \mathbf{a}_2^\top \mathbf{x} = c_2\}$ .

---

**Algorithm 1** REDUCE  $\mathcal{G}_{z,2}$  TO  $\mathcal{MC}_2$ 

---

**Input:**  $(A, c^A, \epsilon_1)$  where  $A \in \mathcal{G}_{z,2}$  is an  $m \times n$  matrix,  $c^A \in \mathbb{R}^n$ ,  $0 < \epsilon_1 < 1$ , and  $\alpha > 0$ .

**Output:**  $(B, c^B, \epsilon_2)$  where  $B \in \mathcal{MC}_2$  is an  $m \times n$  matrix,  $c^B \in \mathbb{R}^n$ , and  $0 < \epsilon_2 < 1$ .

```
1:  $\epsilon_2 \leftarrow \epsilon_1 \left(1 + \frac{1}{\alpha}\right)^{-1/2} \left(1 + \frac{\|c^A\|_2^2 \sigma_{\max}^2(A)}{\alpha+1}\right)^{-1/2}$ 
2:  $\mathcal{X} \leftarrow \{u_1, \dots, u_n, v_1, \dots, v_n\}$  { $\mathcal{MC}_2$  variables and index of new variables}
3: Let  $x$  be the vector of variables corresponding to the set of variables  $\mathcal{X}$ 
4:  $t \leftarrow n + 1$ 
5: Initialize  $\mathcal{A} \leftarrow \emptyset, \mathcal{B} \leftarrow \emptyset$  {Multisets of main and  $\mathcal{MC}_2$  auxiliary equations respectively}
6: for each equation  $1 \leq i \leq m$  in  $A$  do
7:   Let  $\mathcal{I}^{+1} \leftarrow \{j : A_{ij} > 0\}, \mathcal{I}^{-1} \leftarrow \{j : A_{ij} < 0\}$ 
8:   if  $|\mathcal{I}^{+1}| = 1$  and  $|\mathcal{I}^{-1}| = 1$  then
9:     Let the only elements in  $\mathcal{I}^{+1}$  and  $\mathcal{I}^{-1}$  be  $j_+$  and  $j_-$  respectively. {Note
10:     $\mathcal{A} \leftarrow \mathcal{A} \cup \{A_{ij_+} u_t - A_{ij_-} u_{j_-} = c_i\}$ 
11:     $w_i \leftarrow \alpha$ 
12:     $\mathcal{B} \leftarrow \mathcal{B} \cup w_i^{1/2} \cdot \{A_{ij_+} u_{j_+} - A_{ij_-} u_t = 0\}$ 
13:     $\mathcal{X} \leftarrow \mathcal{X} \cup \{u_t, v_t\}$ , update  $x$  accordingly
14:     $t \leftarrow t + 1$ 
15:   else
16:     $\mathcal{A}_i \leftarrow \{A_i u = c_i\}$ 
17:     $\mathcal{B}_i \leftarrow \emptyset$ 
18:    for  $s = -1, +1$  do
19:       $r \leftarrow 0$ 
20:      while  $\mathcal{A}_i$  has strictly more than 1 coefficient with sign  $s$  do
21:        For each  $j$ , let  $\widehat{A}_{ij}$  be the coefficient of  $u_j$  in  $\mathcal{A}_i$ .
22:         $\mathcal{I}_{odd}^s \leftarrow \{j \in \mathcal{I}^s : \lfloor |\widehat{A}_{ij}|/2^r \rfloor \text{ is odd}\}$ 
23:        Pair the indices of  $\mathcal{I}_{odd}^s$  into  $k$  disjoint pairs  $(j_k, l_k)$ 
24:        for each pair of indices  $(j_k, l_k)$  do
25:           $\mathcal{A}_i \leftarrow \mathcal{A}_i + s \cdot 2^r \cdot \{(2u_t - (u_{j_k} + u_{l_k})) = 0\}$ 
26:           $\mathcal{B}_i \leftarrow \mathcal{B}_i \cup -s \cdot 2^r \cdot \mathcal{MC}_2\text{GADGET}(u_{j_k}, u_{l_k}, t)$ 
27:           $\mathcal{X} \leftarrow \mathcal{X} \cup \{u_t, \dots, u_{t+6}, v_t, \dots, v_{t+6}\}$ , update  $x$  accordingly
28:           $t \leftarrow t + 7$ 
29:           $r \leftarrow r + 1$ 
30:        end for
31:      end while
32:    end for
33:     $w_i \leftarrow \alpha |\mathcal{B}_i|$ 
34:     $\mathcal{B} \leftarrow \mathcal{B} \cup w_i^{1/2} \cdot \mathcal{B}_i$ 
35:     $\mathcal{A} \leftarrow \mathcal{A} \cup \mathcal{A}_i$ 
36:  end if
37: end for
38: return  $\epsilon_2$  and  $B, c$  s.t.  $Bx = c$  corresponds to the equations in  $\mathcal{A} \cup \mathcal{B}$ , on the variable set  $\mathcal{X}$ .
```

---

The central object created by Algorithm 1 is the matrix  $B$ , which contains both new equations and new variables. We will superscript the variables with  $A$  to distinguish variables appear in the original equation  $Ax^A = c^A$  from new variables. We will term the new variables as  $x^{\text{aux}}$ , and write a vector solution for the new problem,  $x^B$ , as:

$$x^B = \begin{pmatrix} x^A \\ x^{\text{aux}} \end{pmatrix}. \quad (3.1)$$

Let  $n_A$  be the dimension of  $x^A$ , and  $n_B$  be the dimension of  $x^B$ , respectively. We order the variables so that for an appropriately chosen index  $h$  we have

1.  $x_{1:h}^{\text{aux}}$  corresponds to the  $u$ -coordinates of the auxiliary variables created in  $\mathcal{MC}_2$ -gadgets.
2.  $x_{h+1:n_B-n_A}^{\text{aux}}$  corresponds to the  $v$ -coordinates of the auxiliary variables created in  $\mathcal{MC}_2$ -gadgets.

With this ordering  $x_t^B = u_t$  for  $t \leq n_A + h$ .

Furthermore, we will distinguish the equations in  $B$  into ones formed from manipulating  $A$ , i.e. the equations added to the set  $\mathcal{A}$ , from the auxiliary equations, i.e. the equations added to the set  $\mathcal{B}$ . We use  $W^{1/2} = \text{DIAG}_i(w_i^{1/2})$  to refer to the diagonal matrix of weights  $w_i$  applied to the auxiliary equations  $\mathcal{B}$  in Algorithm 1. In Algorithm 1, a real value  $\alpha > 0$  is set initially and used when computing the weights  $w_i^{1/2}$ . For convenience, throughout most of this section, we will treat  $\alpha$  as an arbitrary constant, and only eventually substitute in its value to complete our main proof.

This leads to the following representation of  $B$  and  $c^B$  which we will use throughout our analysis of the algorithm:

$$B = \begin{pmatrix} \hat{A} \\ W^{1/2} \hat{B} \end{pmatrix}. \quad (3.2)$$

Here the equations of  $\widehat{\mathbf{A}}$  corresponds to  $\mathcal{A}$  in  $\text{REDUCE } \mathcal{G}_{z,2}\text{TO}\mathcal{MC}_2$ , and  $\widehat{\mathbf{B}}$  corresponds to the auxiliary constraints, i.e. equations of  $\mathcal{B}$  in  $\text{REDUCE } \mathcal{G}_{z,2}\text{TO}\mathcal{MC}_2$ . Also, the vector  $\mathbf{c}^{\text{B}}$  created is simply an extension of  $\mathbf{c}^{\text{A}}$ :

$$\mathbf{c}^{\text{B}} = \begin{pmatrix} \mathbf{c}^{\text{A}} \\ \mathbf{0} \end{pmatrix}. \quad (3.3)$$

Finally, as Algorithm 1 creates new equations for each row of  $\mathbf{A}$  independently, we will use  $S_i$  to denote the subset of indices of the rows of  $\widehat{\mathbf{B}}$  that's created from  $\mathbf{A}_{i,:}$ , aka. the the auxiliary constraints generated from the call to  $\mathcal{MC}_2\text{GADGET}$  upon processing  $\mathbf{A}_{i,:}$ . We will also denote the size of these using

$$m_i \stackrel{\text{def}}{=} |S_i|,$$

and use  $\widehat{\mathbf{B}}_i$  to denote these part of  $\widehat{\mathbf{B}}$  that corresponds to these rows. The Gadget routine used in the reduction, and the (trivial) solution mapper are stated below.

---

**Algorithm 2**  $\mathcal{MC}_2\text{GADGET}$

---

**Input:** Scalar variables  $\mathbf{u}_{j_k}, \mathbf{u}_{l_k}$ , integer  $t$ .

**Output:** Set of equations for a multi-commodity gadget that computes the average of  $\mathbf{u}_{j_k}$  and  $\mathbf{u}_{l_k}$ .

1: **return**

$$\begin{aligned} &\{\mathbf{u}_{t+3} - \mathbf{v}_{t+3} - (\mathbf{u}_{t+4} - \mathbf{v}_{t+4}) = 0, \\ &\quad \mathbf{u}_t - \mathbf{u}_{t+3} = 0, \\ &\quad \mathbf{u}_{t+4} - \mathbf{u}_{j_k} = 0, \\ &\quad \mathbf{v}_{t+3} - \mathbf{v}_{t+1} = 0, \\ &\quad \mathbf{v}_{t+2} - \mathbf{v}_{t+4} = 0, \\ &\mathbf{u}_{t+5} - \mathbf{v}_{t+5} - (\mathbf{u}_{t+6} - \mathbf{v}_{t+6}) = 0, \\ &\quad \mathbf{u}_t - \mathbf{u}_{t+5} = 0, \\ &\quad \mathbf{u}_{t+6} - \mathbf{u}_{l_k} = 0, \\ &\quad \mathbf{v}_{t+5} - \mathbf{v}_{t+2} = 0, \\ &\quad \mathbf{v}_{t+1} - \mathbf{v}_{t+6} = 0\} \end{aligned}$$


---

---

**Algorithm 3** MAPSOLN  $\mathcal{MC}_2\text{TO}\mathcal{G}_{z,2}$ 

---

**Input:**  $m \times n$  matrix  $\mathbf{A} \in \mathcal{G}_{z,2}$ ,  $m' \times n'$  matrix  $\mathbf{B} \in \mathcal{MC}_2$ , vector  $\mathbf{c}^A \in \mathbb{R}^m$ , vector  $\mathbf{x}^B \in \mathbb{R}^{n'}$ .

**Output:** Vector  $\mathbf{x}^A \in \mathbb{R}^n$ .

```
1: if  $\mathbf{A}^\top \mathbf{c}^A = \mathbf{0}$  then  
2:   return  $\mathbf{x}^A \leftarrow \mathbf{0}$   
3: else  
4:   return  $\mathbf{x}^A \leftarrow \mathbf{x}_{1:n}^B$   
5: end if
```

---

We upper bound the number of nonzero entries of  $\mathbf{B}$  by the number of nonzero entries of  $\mathbf{A}$ .

**Lemma 3.3.1.**  $\text{nnz}(\mathbf{B}) = O(\text{nnz}(\mathbf{A}) \log \|\mathbf{A}\|_\infty)$  and both dimensions of  $\mathbf{B}$  are  $O(\text{nnz}(\mathbf{A}) \log \|\mathbf{A}\|_\infty)$ .

*Proof.* Since Algorithm 1 constructs new equations for each row of  $\mathbf{A}$  independently, we bound the number of new variables and new equations (that is, the size of the submatrix  $\begin{pmatrix} \widehat{\mathbf{A}}_i \\ \widehat{\mathbf{B}}_i \end{pmatrix}$ ) for each row  $i$  of  $\mathbf{A}$  separately.

Let  $n_i$  be the number of nonzero entries of the  $i$ th row of  $\mathbf{A}$ . We count the number of variables in the equation  $\mathcal{A}_i$ , at each iteration of the while-loop in line 20 of Algorithm 1. Let  $X^{(r)}$  be the subset of variables with nonzero coefficients in  $\mathcal{A}_i$ , at the end of iteration  $r$ . Let  $X_{aux}^{(r)}$  be the subset of  $X^{(r)}$  containing all auxiliary variables created in iteration  $r$ .

Note in each iteration, Algorithm 1 replaces two variables by a new auxiliary variable. It gives that

$$|X_{aux}^{(1)}| \leq \frac{n_i}{2}, \text{ and } |X^{(1)} \setminus X_{aux}^{(1)}| \leq n_i.$$

Since each auxiliary variable in  $X_{aux}^{(1)}$  has coefficient 2, in the 2nd iteration, together with another variable of coefficient 2, it must be replaced by a new auxiliary variable. Thus, at the end of the 2nd iteration, all auxiliary variables  $X_{aux}^{(1)}$  will not appear in the equation  $\mathcal{A}_i$ . This implies that

$$|X_{aux}^{(2)}| \leq \frac{n_i}{2} + \frac{n_i}{4}, \text{ and } |X^{(2)} \setminus X_{aux}^{(2)}| \leq n_i.$$

Similarly, at the end of the  $t^{\text{th}}$  iteration, we have

$$|X_{aux}^{(r)}| \leq \sum_{1 \leq s \leq r} \frac{n_i}{2^s} \leq n_i,$$

and

$$|X^{(r)} \setminus X_{aux}^{(r)}| \leq n_i.$$

Since Algorithm 1 pulls out a factor 2 in each iteration, the total number of iterations is at most  $\log \|\mathbf{A}_i\|_1$ . Since each auxiliary variable in  $\mathcal{A}_i$  during the construction corresponds to  $O(1)$  auxiliary variables and  $O(1)$  equations in  $\widehat{\mathbf{B}}_i$ , the total number of auxiliary variables and equations for row  $i$  of  $\mathbf{A}$  is

$$O(n_i \log \|\mathbf{A}_i\|_1).$$

Therefore, the number of variables and the number of equations in  $\mathbf{B}$  (that is, the both dimensions of  $\mathbf{B}$ ) are

$$O(\text{nnz}(\mathbf{A}) \log \|\mathbf{A}\|_\infty).$$

Since each row of  $\mathbf{B}$  has  $O(1)$  nonzero coefficients, we have

$$\text{nnz}(\mathbf{B}) = O(\text{nnz}(\mathbf{A}) \log \|\mathbf{A}\|_\infty).$$

This completes the proof. □

### *Reduction Between Exact Solvers*

The most important relation between  $\mathbf{A}$  and  $\mathbf{B}$  is given by the following claim.

**Claim 3.3.2** (Reduction between exact solvers). *Fix any  $\mathbf{x}^A \in \mathbb{R}^n$ . Then*

$$\|\mathbf{A}\mathbf{x}^A - \mathbf{c}^A\|_2^2 = \frac{\alpha + 1}{\alpha} \min_{\mathbf{x}^{aux}} \left\| \mathbf{B} \begin{pmatrix} \mathbf{x}^A \\ \mathbf{x}^{aux} \end{pmatrix} - \mathbf{c}^B \right\|_2^2.$$

As a Corollary of Claim 3.3.2, we observe the following:

**Lemma 3.3.3.** *Given LSA problem instance  $(\mathbf{A}, \mathbf{c}^A, \epsilon^A)$  where  $\mathbf{A} \in \mathcal{G}_{z,2}$ , suppose*

$$(\mathbf{B}, \mathbf{c}^B, \epsilon^B) = \text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2(\mathbf{A}, \mathbf{c}^A, \epsilon^A).$$

*Then  $\mathbf{B} \in \mathcal{MC}_2$  and if  $\mathbf{x}^B = \begin{pmatrix} \mathbf{x}^A \\ \mathbf{x}^{aux} \end{pmatrix}$  is a solution to the exact LSA problem  $(\mathbf{B}, \mathbf{c}^B, 0)$ , then  $\mathbf{x}^A$  is a solution to the exact LSA problem  $(\mathbf{A}, \mathbf{c}^A, 0)$ .*

*Proof.* This follows immediately from minimizing over  $\mathbf{x}^A$  on both sides of the equation established by Claim 3.3.2 and then applying and Fact 2.2.2.  $\square$

Before proving Claim 3.3.2 we first note a basic guarantee obtained by  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$ :

$$\mathbf{1}^\top \left( \begin{pmatrix} \hat{\mathbf{A}}_i \\ \hat{\mathbf{B}}_i \end{pmatrix} \mathbf{x}^B - \begin{pmatrix} \mathbf{c}_i^A \\ \mathbf{0} \end{pmatrix} \right) = \mathbf{A}_i \mathbf{x}^A - \mathbf{c}_i^A. \quad (3.4)$$

To verify this guarantee, we consider two cases separately. The first case is when the condition  $|\mathcal{I}^{+1}| = 1$  AND  $|\mathcal{I}^{-1}| = 1$  is true (see Algorithm 1, Line 7). In this case, the main constraint in the output corresponding to row  $i$  is  $\{\mathbf{A}_{ij_+} \mathbf{x}_t^B - \mathbf{A}_{ij_+} \mathbf{x}_{j_-}^B = \mathbf{c}_i^A\}$ , while the auxiliary constraints contain only a single row  $\{\mathbf{A}_{ij_+} \mathbf{x}_{j_+}^B - \mathbf{A}_{ij_+} \mathbf{x}_t^B = 0\}$ , and adding these proves the guarantee for this case. The second case is when the condition in Algorithm 1, Line 7) is false. We consider the case  $s = +1$ . The case  $s = -1$  is proved similarly. Note that we will refer to variables  $j_k$  and  $l_k$  only in the context of a fixed value of  $t$ , which always ensures that they are unambiguously defined. In the case  $s = +1$ , each time



we modify  $\mathcal{A}_i$  by adding  $2^r (2\mathbf{x}_{t+1}^B - (\mathbf{x}_{j_k}^B + \mathbf{x}_{l_k}^B)) = 0$ , we also use  $\mathcal{MC}_2\text{Gadget}$  to create auxiliary constraints that sum up to exactly  $-2^r (2\mathbf{x}_{t+1}^B - (\mathbf{x}_{j_k}^B + \mathbf{x}_{l_k}^B))$ , so adding these together will cancel out the changes.

*Proof of Claim 3.3.2.* For convenience, we also write

1.  $\delta_j \stackrel{\text{def}}{=} \widehat{\mathbf{B}}_j \mathbf{x}^B,$
2.  $\hat{\delta}_i \stackrel{\text{def}}{=} \widehat{\mathbf{A}}_i \mathbf{x}^B - \mathbf{c}_i^A,$
3.  $\epsilon_i \stackrel{\text{def}}{=} \mathbf{A}_i \mathbf{x}^A - \mathbf{c}_i^A.$

Thus

$$\left\| \begin{pmatrix} \widehat{\mathbf{A}}_i \\ \mathbf{W}_i^{1/2} \widehat{\mathbf{B}}_i \end{pmatrix} \mathbf{x}^B - \begin{pmatrix} \mathbf{c}_i^A \\ \mathbf{0} \end{pmatrix} \right\|_2^2 = \hat{\delta}_i^2 + w_i \sum_{j \in S_i} \delta_j^2.$$

Summing over all rows, we get

$$\left\| \mathbf{B} \begin{pmatrix} \mathbf{x}^A \\ \mathbf{x}^{\text{aux}} \end{pmatrix} - \mathbf{c}^B \right\|_2^2 = \sum_i \left( \hat{\delta}_i^2 + w_i \sum_{j \in S_i} \delta_j^2 \right).$$

Similarly, the square of row  $i$  of  $\mathbf{A} \mathbf{x}^A - \mathbf{c}^A$  is  $\epsilon_i^2$  and summing over all rows we get

$$\|\mathbf{A} \mathbf{x}^A - \mathbf{c}^A\|_2^2 = \sum_i \epsilon_i^2.$$

Recalling Equation (3.4), we have

$$\mathbf{1}^\top \left( \begin{pmatrix} \widehat{\mathbf{A}}_i \\ \widehat{\mathbf{B}}_i \end{pmatrix} \mathbf{x}^B - \begin{pmatrix} \mathbf{c}_i^A \\ \mathbf{0} \end{pmatrix} \right) = \mathbf{A}_i \mathbf{x}^A - \mathbf{c}_i^A.$$

Thus

$$\hat{\delta}_i + \sum_{j \in S_i} \delta_j = \epsilon_i.$$

By the Cauchy-Schwarz inequality (applied to two vectors  $\mathbf{a}$  and  $\mathbf{b}$  given by  $\mathbf{a}_1 = \hat{\delta}_i$ ,  $\mathbf{a}_j = \sqrt{w_i}\delta_j$ ,  $\mathbf{b}_1 = 1$ ,  $\mathbf{b}_j = 1/\sqrt{w_i}$ ),

$$\epsilon_i^2 = \left( \hat{\delta}_i + \sum_{j \in S_i} \delta_j \right)^2 \leq \left( \hat{\delta}_i^2 + w_i \sum_{j \in S_i} \delta_j^2 \right) \left( 1 + \frac{m_i}{w_i} \right).$$

the equality holds if and only if

$$\hat{\delta}_i = w_i \delta_j \quad \forall j \in S_i. \quad (3.5)$$

Note that we have ensured that for all  $i$ ,  $w_i = \alpha m_i$ .

By summing over rows we conclude that for every  $\mathbf{x}^A$  and every  $\mathbf{x}^{\text{aux}}$ , we have

$$\|\mathbf{A}\mathbf{x}^A - \mathbf{c}^A\|_2^2 \leq \frac{\alpha + 1}{\alpha} \left\| \mathbf{B} \begin{pmatrix} \mathbf{x}^A \\ \mathbf{x}^{\text{aux}} \end{pmatrix} - \mathbf{c}^B \right\|_2. \quad (3.6)$$

The inequality above will be an equality if Equations (3.5) are satisfied. We now show that for every fixed  $\mathbf{x}^A$ , minimizing over  $\mathbf{x}^{\text{aux}}$  ensures that (3.6) holds with equality.

In particular, we will momentarily prove the following Claim.

**Claim 3.3.4.** *For any fixed  $\mathbf{x}^A$  and its associated  $\epsilon_i$  values, for each row  $i$  of  $\mathbf{A}$ , the linear system*

$$\hat{\mathbf{A}}_i \begin{pmatrix} \mathbf{x}^A \\ \mathbf{x}^{\text{aux}} \end{pmatrix} = \mathbf{c}_i^A + \frac{\alpha}{\alpha + 1} \epsilon_i, \quad (3.7)$$

$$\hat{\mathbf{B}}_j \begin{pmatrix} \mathbf{x}^A \\ \mathbf{x}^{\text{aux}} \end{pmatrix} = \frac{1}{(\alpha + 1)m_i} \epsilon_i, \forall j \in S_i. \quad (3.8)$$

*has a solution (which may not be unique).*

Since every auxiliary variable is associated with only one row  $i$  of  $\mathbf{A}$ , Claim 3.3.4

implies that we can choose  $\mathbf{x}^{\text{aux}}$  s.t. all these linear systems are satisfied at once.

Given such a choice of  $\mathbf{x}^{\text{aux}}$ , we get that Equations (3.5) are satisfied so

$$\|\mathbf{A}\mathbf{x}^{\text{A}} - \mathbf{c}^{\text{A}}\|_2^2 = \frac{\alpha + 1}{\alpha} \left\| \mathbf{B} \begin{pmatrix} \mathbf{x}^{\text{A}} \\ \mathbf{x}^{\text{aux}} \end{pmatrix} - \mathbf{c}^{\text{B}} \right\|_2^2.$$

Given Claim 3.3.4, this completes the proof of Claim 3.3.2.  $\square$

*Proof of Claim 3.3.4.* We focus on the case of the condition  $|\mathcal{I}^{+1}| = 1$  AND  $|\mathcal{I}^{-1}| = 1$  in Algorithm 1, Line 7 being false. The case when then condition is true is very similar, but easier as it deals with a set of just two equations.

We will construct an assignment to all the variables of  $\mathbf{x}^{\text{aux}}$  s.t. Equations (3.7) and (3.8) are satisfied. We start with an assignment  $\mathbf{x}^{\text{A}}$  to the main variables, and we then assign values to auxiliary variables in the order they are created by the algorithm REDUCE  $\mathcal{G}_{z,2}$  TO  $\mathcal{MC}_2$ . Note that we will refer to variables  $j_k$  and  $l_k$  only in the context of a fixed value of  $t$ , which always ensures that they are unambiguously defined. When the algorithm processes pair  $\mathbf{x}_{j_k}^{\text{B}}, \mathbf{x}_{l_k}^{\text{B}} = \mathbf{u}_{j_k}, \mathbf{u}_{l_k}$ , the value of these variables will have been set already, while  $\mathbf{x}_t^{\text{B}} = \mathbf{u}_t$  and the other newly created auxiliary variables have not. Every auxiliary variable is associated with only one row, so we never get multiple assignments to a variable using this procedure.

Recall the constraints created by the  $\mathcal{MC}_2$ Gadget call are

$$\begin{aligned}
& \{\mathbf{u}_{t+3} - \mathbf{v}_{t+3} - (\mathbf{u}_{t+4} - \mathbf{v}_{t+4}) = 0, \\
& \quad \mathbf{u}_t - \mathbf{u}_{t+3} = 0, \\
& \quad \mathbf{u}_{t+4} - \mathbf{u}_{j_k} = 0, \\
& \quad \mathbf{v}_{t+3} - \mathbf{v}_{t+1} = 0, \\
& \quad \mathbf{v}_{t+2} - \mathbf{v}_{t+4} = 0, \\
& \mathbf{u}_{t+5} - \mathbf{v}_{t+5} - (\mathbf{u}_{t+6} - \mathbf{v}_{t+6}) = 0, \\
& \quad \mathbf{u}_t - \mathbf{u}_{t+5} = 0, \\
& \quad \mathbf{u}_{t+6} - \mathbf{u}_{l_k} = 0, \\
& \quad \mathbf{v}_{t+5} - \mathbf{v}_{t+2} = 0, \\
& \quad \mathbf{v}_{t+1} - \mathbf{v}_{t+6} = 0\}
\end{aligned}$$

Let  $\mathbf{z} = (\mathbf{u}_{j_k}, \mathbf{u}_{l_k}, \mathbf{u}_t, \mathbf{u}_{t+1}, \mathbf{u}_{t+2}, \mathbf{u}_{t+3}, \dots, \mathbf{u}_{t+6}, \mathbf{v}_t, \mathbf{v}_{t+1}, \dots, \mathbf{v}_{t+6}) \in \mathbb{R}^{16}$ . Let  $\mathbf{G} \in \mathbb{R}^{10 \times 16}$  be the matrix s.t.  $\mathbf{G}\mathbf{z} = \mathbf{0}$  corresponds to the constraints listed above. Note that all coefficients of  $\mathbf{u}_{t+1}$ ,  $\mathbf{u}_{t+2}$  and  $\mathbf{v}_t$  are zero. We set these three variables to zero.

For some  $\epsilon$ , which we will fix later, we choose  $\mathbf{u}_t$  such that

$$\mathbf{u}_t = \frac{1}{2}(\mathbf{u}_{j_k} + \mathbf{u}_{l_k}) + 5\epsilon.$$

Again, for the same  $\epsilon$ , we fix the following values for  $\mathbf{u}_{t+i}, \mathbf{v}_{t+j}, 3 \leq i \leq 6, 1 \leq j \leq 6$

$$\mathbf{u}_{t+3} = \mathbf{u}_t - \epsilon,$$

$$\mathbf{u}_{t+4} = \mathbf{u}_{j_k} + \epsilon,$$

$$\mathbf{u}_{t+5} = \mathbf{u}_t - \epsilon,$$

$$\mathbf{u}_{t+6} = \mathbf{u}_{l_k} + \epsilon,$$

$$\mathbf{v}_{t+1} = 0,$$

$$\mathbf{v}_{t+2} = 5\epsilon - (\mathbf{u}_t - \mathbf{u}_{j_k}),$$

$$\mathbf{v}_{t+3} = \epsilon,$$

$$\mathbf{v}_{t+4} = 4\epsilon - (\mathbf{u}_t - \mathbf{u}_{j_k}),$$

$$\mathbf{v}_{t+5} = 6\epsilon - (\mathbf{u}_t - \mathbf{u}_{j_k}),$$

$$\mathbf{v}_{t+6} = -\epsilon.$$

This ensures  $\mathbf{G}\mathbf{z} = \epsilon\mathbf{1}$ . Note that for some  $r$ ,  $\widehat{\mathbf{B}}_i \mathbf{x}^B = 2^r \mathbf{G}\mathbf{z} = 2^r \epsilon \mathbf{1}$ , so by choosing  $\epsilon = 2^{-r} \frac{1}{(\alpha+1)m_i} \epsilon_i$ , we can ensure Equations (3.8) are satisfied.

Also

$$\mathbf{c}_i^A + \epsilon_i = \mathbf{1}^\top \begin{pmatrix} \widehat{\mathbf{A}}_i \\ \widehat{\mathbf{B}}_i \end{pmatrix} \mathbf{x}^B = \widehat{\mathbf{A}}_i \mathbf{x}^B + \mathbf{1}^\top \widehat{\mathbf{B}}_i \mathbf{x}^B = \widehat{\mathbf{A}}_i \mathbf{x}^B + \frac{1}{\alpha+1} \epsilon_i.$$

Thus, we have

$$\widehat{\mathbf{A}}_i \mathbf{x}^B = \mathbf{c}_i^A + \frac{\alpha}{\alpha+1} \epsilon_i,$$

which is Equation 3.7.

This completes the proof of the claim. □

**Remark.** The optimal solutions for  $\min_{\mathbf{x}^B} \|\mathbf{B}\mathbf{x}^B - \mathbf{c}^B\|_2$  and  $\min_{\mathbf{x}^A} \|\mathbf{A}\mathbf{x}^A - \mathbf{c}^A\|_2$  have a one-to-one map, however, the optimal values are different:

From the proof of Claim 3.3.2 and Equation (2.2)

$$\|\mathbf{c}^A - \Pi_A \mathbf{c}^A\|_2^2 = \left(1 + \frac{1}{\alpha}\right) \|\mathbf{c}^B - \Pi_B \mathbf{c}^B\|_2^2,$$

where we set the weight  $w_i = \alpha m_i, \forall i$ . Note when  $\alpha \rightarrow \infty$ , the two optimal values approach the same value.

### *Relationship Between Schur Complements*

**Claim 3.3.5.** *Let  $B$  be the output coefficient matrix of Algorithm 1 with input coefficient matrix  $A$  and parameter  $\alpha > 0$ . Then,  $\frac{\alpha}{\alpha+1} A^\top A$  is the Schur complement of  $B^\top B$ .*

*Proof.* Write  $B = \begin{pmatrix} B_1 & B_2 \end{pmatrix}$ , where  $B_1$  is the submatrix corresponding to  $\mathbf{x}^A$  and  $B_2$  is the submatrix corresponding to  $\mathbf{x}^{\text{aux}}$ . Then,

$$B^\top B = \begin{pmatrix} B_1^\top B_1 & B_1^\top B_2 \\ B_2^\top B_1 & B_2^\top B_2 \end{pmatrix}.$$

By Claim 3.3.2, for any fixed  $\mathbf{x}^A$ ,

$$\min_{\mathbf{x}^{\text{aux}}} \mathbf{x}^\top B^\top B \mathbf{x} = \frac{\alpha}{\alpha+1} (\mathbf{x}^A)^\top A^\top A \mathbf{x}^A.$$

By Fact 2.1.2,  $\frac{\alpha}{\alpha+1} A^\top A$  is the Schur complement of  $B^\top B$ . □

### 3.3.1 Approximate solvers

We now show that approximate solvers for  $B$  also translate to approximate solvers for  $A$ .

**Lemma 3.3.6.** *Let  $B\mathbf{x}^B = \mathbf{c}^B$  be the linear system returned by a call to*

*REDUCE  $\mathcal{G}_{z,2}$  TO  $\mathcal{MC}_2(A, \mathbf{c}^A, \epsilon^A)$ , and let  $\epsilon^B$  be the error parameter returned by this call.*

Let  $\mathbf{x}^B$  be a vector such that

$$\|\mathbf{B}\mathbf{x}^B - \Pi_B \mathbf{c}^B\|_2 \leq \epsilon^B \|\Pi_B \mathbf{c}^B\|_2. \quad (3.9)$$

Let  $\tilde{\mathbf{x}}^A$  be the vector returned by a call to  $\text{MAPSOLN } \mathcal{MC}_2\text{TOG}_{z,2}(\mathbf{A}, \mathbf{B}, \mathbf{c}^A, \mathbf{x}^B)$ . Then,

$$\|\mathbf{A}\tilde{\mathbf{x}}^A - \Pi_A \mathbf{c}^A\|_2 \leq \epsilon^A \|\Pi_A \mathbf{c}^A\|_2.$$

*Proof.* We first consider the case  $\mathbf{A}^\top \mathbf{c}^A = \mathbf{0}$ , then Algorithm 3 returns  $\tilde{\mathbf{x}}^A = \mathbf{0}$ , which gives

$$\|\mathbf{A}\tilde{\mathbf{x}}^A - \Pi_A \mathbf{c}^A\|_2 = 0,$$

and completes the proof for this case. When  $\mathbf{A}^\top \mathbf{c}^A \neq \mathbf{0}$ , the Algorithm 3 returns  $\tilde{\mathbf{x}}^A = \mathbf{x}^A$ , where  $\mathbf{x}^B = \begin{pmatrix} \mathbf{x}^A \\ \mathbf{x}^{\text{aux}} \end{pmatrix}$ . Let

$$\mathbf{x}^{B*} \stackrel{\text{def}}{=} \begin{pmatrix} \mathbf{x}^{A*} \\ \mathbf{x}^{\text{aux}*} \end{pmatrix} \in \arg \min_{\mathbf{x}} \|\mathbf{B}\mathbf{x} - \mathbf{c}^B\|_2.$$

By Fact 2.2.2, the condition (3.9) is equivalent to

$$\|\mathbf{x}^B - \mathbf{x}^{B*}\|_{B^\top B} \leq \epsilon^B \|\Pi_B \mathbf{c}^B\|_2.$$

By Claim 3.3.2,  $\mathbf{x}^{A*} \in \arg \min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{c}^A\|_2$ .

We now upper bound the error of  $\|\mathbf{x}^A - \mathbf{x}^{A*}\|_{A^\top A}$ . Note

$$\mathbf{x}^B - \mathbf{x}^{B*} = \begin{pmatrix} \mathbf{x}^A - \mathbf{x}^{A*} \\ \mathbf{x}^{\text{aux}} - \mathbf{x}^{\text{aux}*} \end{pmatrix}.$$

By Claim 2.1.2 and Claim 3.3.5,

$$\|\mathbf{x}^A - \mathbf{x}^{A*}\|_{\mathbf{A}^\top \mathbf{A}}^2 \leq \left(1 + \frac{1}{\alpha}\right) \|\mathbf{x}^B - \mathbf{x}^{B*}\|_{\mathbf{B}^\top \mathbf{B}}^2.$$

Then, we lower bound  $\|\Pi_A \mathbf{c}^A\|_2$ . By Claim 3.3.2,

$$\|\mathbf{c}^A - \Pi_A \mathbf{c}^A\|_2^2 = \left(1 + \frac{1}{\alpha}\right) \|\mathbf{c}^B - \Pi_B \mathbf{c}^B\|_2^2.$$

Since  $\Pi_A \mathbf{c}^A \perp \mathbf{c}^A - \Pi_A \mathbf{c}^A$ , and  $\Pi_B \mathbf{c}^B \perp \mathbf{c}^B - \Pi_B \mathbf{c}^B$ , we have

$$\begin{aligned} \|\Pi_A \mathbf{c}^A\|_2^2 &= \|\mathbf{c}^A\|_2^2 - \|\mathbf{c}^A - \Pi_A \mathbf{c}^A\|_2^2 \\ &= \|\mathbf{c}^B\|_2^2 - \|\mathbf{c}^B - \Pi_B \mathbf{c}^B\|_2^2 - \frac{1}{\alpha + 1} \|\mathbf{c}^A - \Pi_A \mathbf{c}^A\|_2^2 \\ &= \|\Pi_B \mathbf{c}^B\|_2^2 - \frac{1}{\alpha + 1} \|\mathbf{c}^A - \Pi_A \mathbf{c}^A\|_2^2. \end{aligned}$$

Thus,

$$\|\Pi_B \mathbf{c}^B\|_2^2 = \|\Pi_A \mathbf{c}^A\|_2^2 + \frac{1}{\alpha + 1} \|\mathbf{c}^A - \Pi_A \mathbf{c}^A\|_2^2.$$

Since  $\mathbf{A}^\top \mathbf{c}^A \neq \mathbf{0}$ , and because  $\mathbf{A}^\top \mathbf{c}^A$  is integral, by Lemma 2.2.5, we have

$$\|\Pi_A \mathbf{c}^A\|_2^2 \geq \frac{1}{\sigma_{\max}^2(\mathbf{A})}.$$

Thus,

$$\|\Pi_B \mathbf{c}^B\|_2^2 \leq \|\Pi_A \mathbf{c}^A\|_2^2 + \frac{1}{\alpha + 1} \|\mathbf{c}^A\|_2^2 \leq \left(1 + \frac{\|\mathbf{c}^A\|_2^2 \sigma_{\max}^2(\mathbf{A})}{\alpha + 1}\right) \|\Pi_A \mathbf{c}^A\|_2^2.$$

Therefore,

$$\|\mathbf{A} \mathbf{x}^A - \Pi_A \mathbf{c}^A\|_2 \leq \left(1 + \frac{1}{\alpha}\right)^{1/2} \left(1 + \frac{\|\mathbf{c}^A\|_2^2 \sigma_{\max}^2(\mathbf{A})}{\alpha + 1}\right)^{1/2} \epsilon^B \|\Pi_A \mathbf{c}^A\|_2.$$



Finally, by the setting

$$\epsilon^B = \frac{\epsilon^A}{\left(1 + \frac{1}{\alpha}\right)^{1/2} \left(1 + \frac{\|\mathbf{c}^A\|_2^2 \sigma_{\max}^2(\mathbf{A})}{\alpha+1}\right)^{1/2}},$$

we have

$$\|\mathbf{A}\mathbf{x}^A - \mathbf{\Pi}_A \mathbf{c}^A\|_2 \leq \epsilon^A \|\mathbf{\Pi}_A \mathbf{c}^A\|_2.$$

This completes the proof. □

### 3.3.2 Bounding Condition Number of the New Matrix

In this section, we show that the condition number of  $\mathbf{B}$  is upper bounded by the condition number of  $\mathbf{A}$  with a  $\text{poly}(n)$  multiplicative factor.

We first characterize the null space of  $\mathbf{B}$ .

Recall that in Equation (3.1), we write  $\mathbf{x}^B = \begin{pmatrix} \mathbf{x}^A \\ \mathbf{x}^{\text{aux}} \end{pmatrix}$ . In the following, we will employ a *different indexing* of  $\mathbf{x}^{\text{aux}}$  than the one defined at the beginning of Section 3.3. We also reorder the columns of  $\mathbf{B}$  so that  $\mathbf{B}\mathbf{x}^B = \mathbf{0}$  represents the same equations as before. For appropriately chosen indices  $g_1$  and  $g_2$ , we define

1.  $\mathbf{x}_{1:g_1}^{\text{aux}}$  corresponds to the  $\mathbf{u}$ -coordinates of the auxiliary variables created in  $\mathcal{MC}_2$ -gadgets *whose coefficients are nonzero*.
2.  $\mathbf{x}_{g_1+1:g_2}^{\text{aux}}$  corresponds to the  $\mathbf{v}$ -coordinates of the auxiliary variables created in  $\mathcal{MC}_2$ -gadgets *whose coefficients are nonzero*.
3.  $\mathbf{x}_{g_2+1:n_B-n_A}^{\text{aux}}$  corresponds to the coordinates of the auxiliary variables created in  $\mathcal{MC}_2$ -gadgets *whose coefficients are zero*.

Using this ordering, for  $0 \leq i \leq (g_2 - g_1)/6$ ,  $\mathbf{x}_{g_1+6i+1:g_1+6i+6}^{\text{aux}}$  corresponds to the  $\mathbf{v}$ -coordinates of the auxiliary variables with non-zero coefficients in a single  $\mathcal{MC}_2$ -gadget.

Given any fixed  $\mathbf{x}^A \in \text{null}(\mathbf{A})$ , we extend  $\mathbf{x}^A$  to a vector in dimension  $n_B$ :

$$\mathbf{p}(\mathbf{x}^A) \stackrel{\text{def}}{=} \begin{pmatrix} \mathbf{x}^A \\ \mathbf{x}^{\text{aux}} \end{pmatrix}. \quad (3.10)$$

We assign the values of the auxiliary variables of  $\mathbf{p}(\mathbf{x}^A)$  in the order that they are created in Algorithm 1. In a  $\mathcal{MC}_2$ -gadget, suppose the values of variables  $\mathbf{u}_{j_k}$  and  $\mathbf{u}_{l_k}$  have already been assigned. Let  $\mathbf{u}_t, \mathbf{u}_{t+i}, \mathbf{v}_{t+j}, 3 \leq i \leq 6, 1 \leq j \leq 6$  be the auxiliary variables created in this gadget. We assign values as follows,

$$\begin{aligned} \mathbf{u}_t &= (\mathbf{u}_{j_k} + \mathbf{u}_{l_k})/2, \\ \mathbf{u}_{t+3} &= \mathbf{u}_t, \\ \mathbf{u}_{t+4} &= \mathbf{u}_{j_k}, \\ \mathbf{u}_{t+5} &= \mathbf{u}_t, \\ \mathbf{u}_{t+6} &= \mathbf{u}_{l_k}, \\ \mathbf{v}_{t+1} &= 0, \\ \mathbf{v}_{t+2} &= \mathbf{u}_t - \mathbf{u}_{j_k}, \\ \mathbf{v}_{t+3} &= 0, \\ \mathbf{v}_{t+4} &= \mathbf{u}_t - \mathbf{u}_{j_k}, \\ \mathbf{v}_{t+5} &= \mathbf{u}_t - \mathbf{u}_{j_k}, \\ \mathbf{v}_{t+6} &= 0. \end{aligned}$$

This gives the first  $g_2$  entries of  $\mathbf{x}^{\text{aux}}$ , and we set all the rest of the entries to be 0.

Let  $\mathbf{e}_i \in \mathbb{R}^{(g_2-g_1)/6}$  be the  $i$ th standard basis vector and  $\mathbf{1}$  be the all-one vector in 6 dimensions. We define  $\mathbf{p}^i \in \mathbb{R}^{n_B}$  to be a vector whose nonzero entries are given by

$$\mathbf{p}_{n_A+g_1+1:n_A+g_2}^i = \mathbf{e}_i \otimes \mathbf{1}. \quad (3.11)$$

Let  $\mathbf{e}_j \in \mathbb{R}^{n_B - n_A - g_2}$  is the  $j$ th standard basis vector. We define  $\mathbf{q}^j \in \mathbb{R}^{n_B}$  to be a vector whose nonzero entries are given by

$$\mathbf{q}_{n_A + g_2 + 1:n_B}^j = \mathbf{e}_j. \quad (3.12)$$

**Lemma 3.3.7.**  $\text{null}(\mathbf{B}) = \text{Span}(\mathbf{p}(\mathbf{x}^A), \mathbf{p}^i, \mathbf{q}^j : \mathbf{x}^A \in \text{null}(\mathbf{A}), 1 \leq i \leq \frac{g_2 - g_1}{6}, 1 \leq j \leq n_B - g_2).$

*Proof.* Let

$$\mathbf{S} \stackrel{\text{def}}{=} \text{Span}\left(\mathbf{p}(\mathbf{x}^A), \mathbf{p}^i, \mathbf{q}^j : \mathbf{x}^A \in \text{null}(\mathbf{A}), 1 \leq i \leq \frac{g_2 - g_1}{6}, 1 \leq j \leq n_B - g_2\right).$$

According to definitions of the vectors, we can check that  $\mathbf{S} \subseteq \text{null}(\mathbf{B})$ .

It remains to show that  $\mathbf{S} \supseteq \text{null}(\mathbf{B})$ . Let  $\mathbf{x} \in \text{null}(\mathbf{B})$ . By Claim 3.3.2 with  $\mathbf{c}^A = \mathbf{0}$ , we have

$$\mathbf{A}\mathbf{x}^A = \mathbf{0},$$

that is,  $\mathbf{x}^A \in \text{null}(\mathbf{A})$ . According to the  $\mathcal{MC}_2$ -gadget constraints in Algorithm 2, we have

$$\mathbf{v}_{t+1} = \mathbf{v}_{t+3} = \mathbf{v}_{t+6} = \gamma, \text{ and } \mathbf{v}_{t+2} = \mathbf{v}_{t+4} = \mathbf{v}_{t+5} = \theta,$$

where  $\gamma - \theta = (\mathbf{u}_{j_k} - \mathbf{u}_{j_l})/2$ . Besides, since all entries in  $\mathbf{x}_{n_A + g_2 + 1:n_B}$  have zero coefficients, they are free to choose. Thus,  $\mathbf{x} \in \mathbf{S}$ , that is,  $\mathbf{S} \supseteq \text{null}(\mathbf{B})$ .

This completes the proof. □

We upper bound the largest singular value of  $\mathbf{B}$  in the following claim.

**Lemma 3.3.8.**  $\lambda_{\max}(\mathbf{B}^\top \mathbf{B}) = O(\alpha \text{nnz}(\mathbf{A})^2 \|\mathbf{A}\|_\infty^2 \log^2 \|\mathbf{A}\|_\infty).$

*Proof.* By the Courant-Fischer Theorem,

$$\lambda_{\max}(\mathbf{B}^\top \mathbf{B}) = \max_{\mathbf{x}} \frac{\mathbf{x}^\top \mathbf{B}^\top \mathbf{B} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}.$$

We write  $\mathbf{x}$  as  $(\mathbf{u}_1, \mathbf{v}_2, \dots, \mathbf{u}_{n_B}, \mathbf{v}_{n_B})$ , so that each  $(\mathbf{u}_i, \mathbf{v}_i)$  corresponds to the two coordinates of vertex  $i$  in the graph. Expanding the right hand side, we get

$$\frac{\sum_{(i,j) \in E_1} w_{(i,j)} (\mathbf{u}_i - \mathbf{u}_j)^2 + \sum_{(i,j) \in E_2} w_{(i,j)} (\mathbf{v}_i - \mathbf{v}_j)^2 + \sum_{(i,j) \in E_{1+2}} w_{(i,j)} (\mathbf{u}_i - \mathbf{v}_i - (\mathbf{u}_j - \mathbf{v}_j))^2}{\sum_i \mathbf{u}_i^2 + \mathbf{v}_i^2}$$

where  $w_{(i,j)}$ 's are the edge weights. Let

$$w_{\max} = \max_{(i,j) \in E_1 \cup E_2 \cup E_{1+2}} w_{(i,j)}.$$

According to our construction of  $\mathbf{B}$  in Algorithm 1,

$$w_{\max} \leq \alpha \text{nnz}(\mathbf{B}) \|\mathbf{A}\|_{\infty}^2. \quad (3.13)$$

Thus,  $\lambda_{\max}(\mathbf{B}^{\top} \mathbf{B})$  is at most

$$w_{\max} \max_{\mathbf{u}} \frac{\sum_{(i,j) \in E_1} (\mathbf{u}_i - \mathbf{u}_j)^2 + \sum_{(i,j) \in E_2} (\mathbf{v}_i - \mathbf{v}_j)^2 + \sum_{(i,j) \in E_{1+2}} (\mathbf{u}_i - \mathbf{v}_i - (\mathbf{u}_j - \mathbf{v}_j))^2}{\sum_i \mathbf{u}_i^2 + \mathbf{v}_i^2}.$$

We upper bound each term of the numerator. Let  $d_1, d_2, d_{1+2}$  be the maximum vertex degree of the graphs  $G_1, G_2, G_{1+2}$ , respectively. By Cauchy-Schwarz inequality,

$$\sum_{(i,j) \in E_1} (\mathbf{u}_i - \mathbf{u}_j)^2 \leq 2 \sum_{(i,j) \in E_1} \mathbf{u}_i^2 + \mathbf{u}_j^2 \leq 2d_1 \left( \sum_i \mathbf{u}_i^2 + \mathbf{v}_i^2 \right),$$

$$\sum_{(i,j) \in E_2} (\mathbf{v}_i - \mathbf{v}_j)^2 \leq 2d_2 \left( \sum_i \mathbf{u}_i^2 + \mathbf{v}_i^2 \right),$$

and

$$\begin{aligned} \sum_{(i,j) \in E_{1+2}} (\mathbf{u}_i - \mathbf{v}_i - (\mathbf{u}_j - \mathbf{v}_j))^2 &\leq 4 \sum_{(i,j) \in E_{1+2}} \mathbf{u}_i^2 + \mathbf{v}_i^2 + \mathbf{u}_j^2 + \mathbf{v}_j^2 \\ &\leq 4d_{1+2} \left( \sum_i \mathbf{u}_i^2 + \mathbf{v}_i^2 \right). \end{aligned}$$

Plugging the above inequalities to the expansion, we have

$$\begin{aligned} \lambda_{\max}(\mathbf{B}^\top \mathbf{B}) &\leq w_{\max} \cdot 8 \max\{d_1, d_2, d_{1+2}\} \\ &\leq 8w_{\max} \text{nnz}(\mathbf{B}). \end{aligned}$$

By the upper bound of  $\text{nnz}(\mathbf{B})$  in Lemma 3.3.1 and the upper bound of  $w_{\max}$  in Equation (3.13), we have

$$\lambda_{\max}(\mathbf{B}^\top \mathbf{B}) \leq O(\alpha \text{nnz}(\mathbf{A})^2 \|\mathbf{A}\|_\infty^2 \log^2 \|\mathbf{A}\|_\infty).$$

This completes the proof. □

Recall that, for a vector  $\mathbf{x} \in \mathbb{R}^{n_B}$ ,  $\mathbf{x}_{n_A+1:n_A+g_1}$  corresponds to the auxiliary  $\mathbf{u}$ -variables with non-zero coefficients in  $\mathcal{MC}_2$ -gadgets,  $\mathbf{x}_{n_A+1+g_1+6i:n_A+6+g_1+6i}$  for  $0 \leq i \leq (g_2 - g_1)/6$  corresponds to the auxiliary  $\mathbf{v}$ -variables with non-zero coefficients in a single  $\mathcal{MC}_2$ -gadget, and  $\mathbf{x}_{n+g_2+1:n_B}$  corresponds to the variables with zero coefficient.

**Lemma 3.3.9.** *Let  $0 \leq \epsilon \leq (100n_B \log \|\mathbf{A}\|_\infty)^{-1}$ . Let  $\mathbf{x}^B \in \mathbb{R}^{n_B}$  satisfying*

1.  $\mathbf{x}_{n_A+g_2+1:n_B}^B = \mathbf{0}$ ,
2.  $\mathbf{x}_{n_A+g_1+6i+1}^B = 0, \forall 0 \leq i \leq (g_2 - g_1)/6$ , and
3.  $\|\mathbf{B}\mathbf{x}^B\|_2 \leq \epsilon \|\mathbf{x}^B\|_2$ .

Then,

$$\|\mathbf{A}\mathbf{x}^A\|_2 \leq 4 \left(1 + \frac{1}{\alpha}\right)^{1/2} n_B^{1/2} \epsilon \|\mathbf{x}^A\|_2.$$

*Proof.* We first show that under the conditions of the Lemma,

$$\|\mathbf{x}^A\|_\infty \geq \frac{1}{4} \|\mathbf{x}^B\|_\infty. \quad (3.14)$$

Let  $\delta \stackrel{\text{def}}{=} \epsilon \|\mathbf{x}^B\|_2$ . By the 3rd condition, each entry of  $\mathbf{B}\mathbf{x}^B$  has absolute value at most  $\delta$ .

We first show that all  $\mathbf{u}$ -variables in  $\mathbf{x}^{\text{aux}}$  cannot be large. According to Algorithm 1, for each row of  $\mathbf{B}$ , all nonzero coefficients have same absolute value, which is at least 1. Based on this fact and the  $\mathcal{MC}_2$ -gadget constructed in Algorithm 2, we have

$$\begin{aligned} \frac{\mathbf{u}_{j_k} + \mathbf{u}_{l_k}}{2} - 5\delta &\leq \mathbf{u}_t \leq \frac{\mathbf{u}_{j_k} + \mathbf{u}_{l_k}}{2} + 5\delta \\ \frac{\mathbf{u}_{j_k} + \mathbf{u}_{l_k}}{2} - 6\delta &\leq \mathbf{u}_{t+3} \leq \frac{\mathbf{u}_{j_k} + \mathbf{u}_{l_k}}{2} + 6\delta \\ \mathbf{u}_{j_k} - \delta &\leq \mathbf{u}_{t+4} \leq \mathbf{u}_{j_k} + \delta \\ \frac{\mathbf{u}_{j_k} + \mathbf{u}_{l_k}}{2} - 6\delta &\leq \mathbf{u}_{t+5} \leq \frac{\mathbf{u}_{j_k} + \mathbf{u}_{l_k}}{2} + 6\delta \\ \mathbf{u}_{l_k} - \delta &\leq \mathbf{u}_{t+6} \leq \mathbf{u}_{l_k} + \delta \end{aligned}$$

where  $\mathbf{u}_{j_k}, \mathbf{u}_{l_k}$  being paired-and-replaced, and all others are entries of  $\mathbf{x}^{\text{aux}}$ . By the triangle inequality,

$$|\mathbf{u}_t|, |\mathbf{u}_{t+k}| \leq \frac{1}{2} (|\mathbf{u}_{j_k}| + |\mathbf{u}_{l_k}|) + 6\delta, \quad \forall 3 \leq k \leq 6.$$

Note that the sum of coefficients on both sides are equal. We can repeat this type of substitution on the right hand side until  $\mathbf{u}_{j_k}$  and  $\mathbf{u}_{l_k}$  are variables of  $\mathbf{x}^A$ . At the  $i$ th iteration of Algorithm 1 line 20, we have

$$|\mathbf{u}_t|, |\mathbf{u}_{t+k}| \leq \sum_j \alpha_j |\mathbf{x}_j^A| + 6r\delta, \quad \forall 3 \leq k \leq 6,$$

where  $\alpha_j \geq 0$  and  $\sum_j \alpha_j = 1$ . By the Hölder inequality,

$$|\mathbf{u}_t|, |\mathbf{u}_{t+k}| \leq \|\mathbf{x}^A\|_\infty + 6r\delta, \quad \forall 3 \leq k \leq 6. \quad (3.15)$$

We then argue that all  $\mathbf{v}$ -variables in  $\mathbf{x}^{\text{aux}}$  cannot be large. Note that by the 2nd condition, in the  $\mathcal{MC}_2$ -gadget, we have  $\mathbf{v}_{t+1} = 0$ . According to the equations in Algorithm 2, we have

$$\begin{aligned} \mathbf{u}_{t+3} - \mathbf{u}_{t+4} - 3\delta &\leq \mathbf{v}_{t+2} \leq \mathbf{u}_{t+3} - \mathbf{u}_{t+4} + 3\delta \\ -\delta &\leq \mathbf{v}_{t+3} \leq \delta \\ \mathbf{u}_{t+3} - \mathbf{u}_{t+4} - 2\delta &\leq \mathbf{v}_{t+4} \leq \mathbf{u}_{t+3} - \mathbf{u}_{t+4} + 2\delta \\ \mathbf{u}_{t+3} - \mathbf{u}_{t+4} - 4\delta &\leq \mathbf{v}_{t+5} \leq \mathbf{u}_{t+3} - \mathbf{u}_{t+4} + 4\delta \\ -\delta &\leq \mathbf{v}_{t+6} \leq \delta \end{aligned}$$

By the triangle inequality,

$$|\mathbf{v}_{t+k}| \leq |\mathbf{u}_{t+3}| + |\mathbf{u}_{t+4}| + 4\delta, \quad \forall 2 \leq k \leq 6.$$

By (3.15), at the  $r$ th iteration of Algorithm 1 line 20,

$$|\mathbf{v}_{t+k}| \leq 2\|\mathbf{x}^A\|_\infty + 12r\delta, \quad \forall 2 \leq k \leq 6.$$

Since there are at most  $\log \|\mathbf{A}\|_\infty$  iterations, the above inequality together with Equation (3.15) implies

$$\|\mathbf{x}^{\text{aux}}\|_\infty \leq 2\|\mathbf{x}^A\|_\infty + 12\delta \log \|\mathbf{A}\|_\infty.$$

Adding  $\|\mathbf{x}^A\|_\infty$  on both sides and substituting  $\delta = \epsilon \|\mathbf{x}^B\|_2$  gives

$$\begin{aligned}\|\mathbf{x}^B\|_\infty &\leq 3\|\mathbf{x}^A\|_\infty + 12\epsilon \|\mathbf{x}^B\|_2 \log \|\mathbf{A}\|_\infty \\ &\leq 3\|\mathbf{x}^A\|_\infty + 12\epsilon \sqrt{n_B} \|\mathbf{x}^B\|_\infty \log \|\mathbf{A}\|_\infty.\end{aligned}$$

Given  $\epsilon \leq (100\sqrt{n_B} \log \|\mathbf{A}\|_\infty)^{-1}$ , we have Equation (3.14). This further implies

$$\|\mathbf{x}^B\|_2 \leq 4\sqrt{n_B} \|\mathbf{x}^A\|_2. \quad (3.16)$$

By Claim 3.3.2 with  $\mathbf{c}^B = \mathbf{0}$ , we have

$$\|\mathbf{A}\mathbf{x}^A\|_2 \leq \sqrt{\frac{\alpha+1}{\alpha}} \|\mathbf{B}\mathbf{x}^B\|_2.$$

By the 3rd condition and the bound (3.16), we have

$$\|\mathbf{A}\mathbf{x}^A\|_2 \leq 4 \left(1 + \frac{1}{\alpha}\right)^{1/2} n_B^{1/2} \epsilon \|\mathbf{x}^A\|_2.$$

This completes the proof. □

We use the above lemma to lower bound the smallest nonzero eigenvalue of  $\mathbf{B}^\top \mathbf{B}$ .

**Lemma 3.3.10.**  $\lambda_{\min}(\mathbf{B}^\top \mathbf{B}) = \Omega\left(\frac{\min\{1, \lambda_{\min}(\mathbf{A}^\top \mathbf{A})\}}{n_B^2}\right).$

*Proof.* Let

$$\delta \stackrel{\text{def}}{=} \min\{1, \lambda_{\min}(\mathbf{A}^\top \mathbf{A})\}, \quad (3.17)$$

and

$$\epsilon \stackrel{\text{def}}{=} \frac{\delta}{1616(1 + \alpha^{-1})n_B^2}. \quad (3.18)$$



The goal is to prove that  $\lambda_{\min}(\mathbf{B}^\top \mathbf{B}) \geq \epsilon$ . Assume by contradiction, there exists an  $\mathbf{x} \perp \text{null}(\mathbf{B})$  such that

$$(\mathbf{x}^{\text{B}})^\top \mathbf{B}^\top \mathbf{B} \mathbf{x}^{\text{B}} < \epsilon (\mathbf{x}^{\text{B}})^\top \mathbf{x}^{\text{B}}. \quad (3.19)$$

We show a contradiction by case analysis according to whether  $\mathbf{x}^{\text{A}}$  is orthogonal to the null space of  $\mathbf{A}$ .

**Case 1:** Suppose  $\mathbf{x}^{\text{A}} \perp \text{null}(\mathbf{A})$ .

Recall that  $\mathbf{x}_{n_A+1:n_A+g_1}^{\text{B}}$  corresponds to the auxiliary  $\mathbf{u}$ -variables in the constraints,  $\mathbf{x}_{n_A+g_1+1:n_A+g_2}^{\text{B}}$  corresponds to the auxiliary  $\mathbf{v}$ -variables in the constraints, and  $\mathbf{x}_{n_A+g_2+1:n_B}^{\text{B}}$  corresponds to the variables with zero coefficient. By Lemma 3.3.7, we have  $\mathbf{x}^{\text{B}} \perp \mathbf{q}^j, \forall 1 \leq j \leq n_A - g_2$ , where  $\mathbf{q}^j$  is defined in Equation (3.12). Thus,  $\mathbf{x}_{n_A+g_2+1:n_B} = \mathbf{0}$ .

Let  $\mathbf{r} \stackrel{\text{def}}{=} \sum_{0 \leq i \leq (g_2 - g_1)/6} \mathbf{x}_{n_A+g_1+6i+1}^{\text{B}} \mathbf{p}^i$ , where  $\mathbf{p}^i$  is defined in Equation (3.11). By Lemma 3.3.7,  $\mathbf{r} \in \text{null}(\mathbf{B})$ . Let  $\mathbf{y}^{\text{B}} \stackrel{\text{def}}{=} \mathbf{x}^{\text{B}} - \mathbf{r}$ . Similarly, we write

$$\mathbf{y}^{\text{B}} = \begin{pmatrix} \mathbf{y}^{\text{A}} \\ \mathbf{y}^{\text{aux}} \end{pmatrix},$$

where  $\mathbf{y}^{\text{A}}$  corresponds to original variables and  $\mathbf{y}^{\text{aux}}$  corresponds to auxiliary variables. Note  $\mathbf{y}^{\text{A}} = \mathbf{x}^{\text{A}}$ . Since  $\mathbf{r} \in \text{null}(\mathbf{B})$ , we have

$$(\mathbf{y}^{\text{B}})^\top \mathbf{B}^\top \mathbf{B} \mathbf{y}^{\text{B}} = (\mathbf{x}^{\text{B}})^\top \mathbf{B}^\top \mathbf{B} \mathbf{x}^{\text{B}}.$$

Since  $\mathbf{x}^{\text{B}} \perp \mathbf{r}$ , we have

$$(\mathbf{y}^{\text{B}})^\top \mathbf{y}^{\text{B}} = (\mathbf{x}^{\text{B}})^\top \mathbf{x}^{\text{B}} + \mathbf{r}^\top \mathbf{r} \geq (\mathbf{x}^{\text{B}})^\top \mathbf{x}^{\text{B}}.$$

By the assumption in Equation (3.19), we have

$$(\mathbf{y}^B)^\top \mathbf{B}^\top \mathbf{B} \mathbf{y}^B < \epsilon (\mathbf{y}^B)^\top \mathbf{y}^B.$$

By our definition of  $\mathbf{y}^B$ ,  $\mathbf{y}^B$  satisfies the conditions of Lemma 3.3.9. By Lemma 3.3.9 and the definition of  $\epsilon$  in Equation (3.18), we have

$$(\mathbf{x}^A)^\top \mathbf{A}^\top \mathbf{A} \mathbf{x}^A < \delta (\mathbf{x}^A)^\top \mathbf{x}^A.$$

Since  $\mathbf{x}^A \perp \text{null}(\mathbf{A})$  and the definition of  $\delta$  in Equation (3.17), we have

$$\lambda_{\min}(\mathbf{A}^\top \mathbf{A}) < \delta = \min\{1, \lambda_{\min}(\mathbf{A}^\top \mathbf{A})\},$$

which is a contradiction.

**Case 2:** Suppose  $\mathbf{x}^A \in \text{null}(\mathbf{A})$ .

Recall that  $\mathbf{p}(\mathbf{x}^A)$  is the extension of vector  $\mathbf{x}^A$  defined in Equation (3.10). Let  $\mathbf{r} = \sum_{0 \leq i \leq (g_2 - g_1)/6} \alpha_i \mathbf{p}^i$  be a vector such that

$$\mathbf{y}^B \stackrel{\text{def}}{=} \mathbf{x}^B - \mathbf{p}(\mathbf{x}^A) - \mathbf{r}$$

satisfying  $\mathbf{y}_{n_A + g_1 + 6i + 1}^B = \mathbf{0}, \forall 0 \leq i < (g_2 - g_1)/6$ . By this definition, we have

$$\mathbf{y}^A = \mathbf{0}. \tag{3.20}$$

Since  $\mathbf{p}(\mathbf{x}^A) + \mathbf{r}$  is in the null space of  $\mathbf{B}$ , by the assumption in Equation (3.19),

$$(\mathbf{y}^B)^\top \mathbf{B}^\top \mathbf{B} \mathbf{y}^B < \epsilon (\mathbf{y}^B)^\top \mathbf{y}^B.$$

Note  $\mathbf{y}^B$  satisfies the conditions of Lemma 3.3.9. By Equations (3.14) and (3.20), we have

$$\|\mathbf{y}^B\|_\infty = 0.$$

This implies that  $\mathbf{x}^B = \mathbf{p}(\mathbf{x}^A) + \mathbf{r}$  is in the null space of  $\mathbf{B}$ , which contradicts that  $\mathbf{x}^B \perp \text{null}(\mathbf{B})$ .

**Case 3:** Suppose  $\mathbf{x}^A = \tilde{\mathbf{z}} + \hat{\mathbf{z}}$ , where  $\tilde{\mathbf{z}} \neq \mathbf{0}$  is in  $\text{null}(\mathbf{A})$  and  $\hat{\mathbf{z}} \neq \mathbf{0}$  is orthogonal to  $\text{null}(\mathbf{A})$ .

Similarly, let  $\mathbf{r}_1 = \sum_i \alpha_i \mathbf{p}^i$  such that  $\mathbf{y}^B \stackrel{\text{def}}{=} \mathbf{x}^B - \mathbf{r}_1$  satisfying  $\mathbf{y}_{n_A+g_1+6i+1}^B = 0, \forall i$ .

By this definition, we have

$$\mathbf{y}^A = \mathbf{x}^A.$$

Then by the assumption in Equation (3.19),

$$(\mathbf{y}^B)^\top \mathbf{B}^\top \mathbf{B} \mathbf{y}^B < \epsilon (\mathbf{y}^B)^\top \mathbf{y}^B.$$

$\mathbf{y}^B$  satisfies the conditions of Lemma 3.3.9. By Lemma 3.3.9, we have

$$(\mathbf{y}^A)^\top \mathbf{A}^\top \mathbf{A} \mathbf{y}^A \leq 16 \left(1 + \frac{1}{\alpha}\right) n_B \epsilon (\mathbf{y}^A)^\top \mathbf{y}^A.$$

Since  $\mathbf{y}^A = \tilde{\mathbf{z}} + \hat{\mathbf{z}}$  with  $\tilde{\mathbf{z}} \in \text{null}(\mathbf{A})$  and  $\hat{\mathbf{z}} \perp \text{null}(\mathbf{A})$ , we have

$$\hat{\mathbf{z}}^\top \mathbf{A}^\top \mathbf{A} \hat{\mathbf{z}} \leq 16 \left(1 + \frac{1}{\alpha}\right) n_B \epsilon (\tilde{\mathbf{z}}^\top \tilde{\mathbf{z}} + \hat{\mathbf{z}}^\top \hat{\mathbf{z}}).$$

If  $\hat{\mathbf{z}}^\top \hat{\mathbf{z}} > \tilde{\mathbf{z}}^\top \tilde{\mathbf{z}} / (100n_B)$ , then we have a contradiction with Equation (3.17) and we have done. Otherwise, we have

$$\|\hat{\mathbf{z}}\|_2^2 \leq \frac{\|\tilde{\mathbf{z}}\|_2^2}{100n_B}.$$

Let  $\mathbf{r}_2 = \sum_i \beta_i \mathbf{p}^i$  such that

$$\mathbf{z}^B \stackrel{\text{def}}{=} \mathbf{y}^B - \mathbf{p}(\tilde{\mathbf{z}}) - \mathbf{r}_2$$

satisfies  $\mathbf{z}_{n_A+g_1+6i+1}^B = 0, \forall i$ . By the assumption in Equation (3.19),

$$(\mathbf{z}^B)^\top \mathbf{B}^\top \mathbf{B} \mathbf{z}^B < \epsilon (\mathbf{z}^B)^\top \mathbf{z}^B.$$

Vector  $\mathbf{z}^B$  satisfies the conditions of Lemma 3.3.9. By Equation (3.14), we have

$$\|\mathbf{z}^B\|_2 \leq 4\sqrt{n_B} \|\hat{\mathbf{z}}\|_2.$$

Since  $\hat{\mathbf{z}}^\top \hat{\mathbf{z}} \leq \tilde{\mathbf{z}}^\top \tilde{\mathbf{z}} / (100n_B)$ , we have

$$\|\mathbf{z}^B\|_2 \leq \|\tilde{\mathbf{z}}\|_2.$$

On the other hand,

$$\begin{aligned} \|\mathbf{z}^B\|_2^2 &= \|\mathbf{x}^B\|_2^2 + \|\mathbf{p}(\tilde{\mathbf{z}}) + \mathbf{r}_1 + \mathbf{r}_2\|_2^2 \\ &> \|\mathbf{p}(\tilde{\mathbf{z}}) + \mathbf{r}_1 + \mathbf{r}_2\|_2^2 \\ &\geq \|\mathbf{x}^A\|_2^2 \\ &= \|\tilde{\mathbf{z}}\|_2^2 + \|\hat{\mathbf{z}}\|_2^2 \\ &> \|\tilde{\mathbf{z}}\|_2^2. \end{aligned}$$

The first equality is due to  $\mathbf{x}^B \perp (\mathbf{p}(\tilde{\mathbf{z}}) + \mathbf{r}_1 + \mathbf{r}_2)$ , the third inequality is due to the  $\mathbf{x}^A$  part of  $\mathbf{r}_1 + \mathbf{r}_2$  is  $\mathbf{0}$ , and the fourth equality is due to that  $\tilde{\mathbf{z}} \perp \hat{\mathbf{z}}$ . Thus, we get a contradiction.

This completes the proof. □

Lemma 3.3.8 and Lemma 3.3.10 immediately imply the following lemma.

**Lemma 3.3.11.**  $\kappa(\mathbf{B}) = O\left(\frac{\text{nnz}(\mathbf{A})^2 \|\mathbf{A}\|_\infty \log^2 \|\mathbf{A}\|_\infty}{\min\{1, \sigma_{\min}(\mathbf{A})\}}\right).$

### 3.3.3 Putting it All Together

In this section, we prove Lemma 3.2.5.

*Proof of Lemma 3.2.5.* We set  $\alpha = 1$  in Algorithm 1.

By Lemma 3.3.1 and Claim 3.1.2, we have

$$\text{nnz}(\mathbf{B}) = O(s \log(sU)).$$

According to our reduction in Algorithm 1, the largest entry and the smallest nonzero entry of the right hand side vector does not change. Besides, all nonzero entries of  $\mathbf{B}$  have absolute value at least 1. The largest entry of  $\mathbf{B}$  is upper bounded by Equation (3.13)

$$\sqrt{w_{\max}} = O\left(\sqrt{\text{nnz}(\mathbf{B})} \|\mathbf{A}\|_{\infty}\right).$$

By Claim 3.1.2, we have

$$\sqrt{w_{\max}} = O\left(s^{3/2} U \log^{1/2}(sU)\right).$$

By Lemma 3.3.11 and Lemma 3.1.2, 3.1.2, we have

$$\kappa(\mathbf{B}) = O\left(\frac{\text{nnz}(\mathbf{A})^2 \|\mathbf{A}\|_{\infty} \log^2 \|\mathbf{A}\|_{\infty}}{\min\{1, \sigma_{\min}(\mathbf{A})\}}\right) = O(s^4 U^2 K \log^2(sU)).$$

By Lemma 3.3.6, we have

$$(\epsilon^{\mathbf{B}})^{-1} = (\epsilon^{\mathbf{A}})^{-1} O\left(1 + \|\mathbf{c}^{\mathbf{A}}\|_2 \sigma_{\max}(\mathbf{A})\right) = (\epsilon^{\mathbf{A}})^{-1} O(sU^2).$$

This completes the proof. □

### 3.4 $\mathcal{MC}_2$ Efficiently Reducible to $\mathcal{MC}_2^{>0}$

The matrices generated by interior point methods (see Section 3.8 for details) are more restrictive than  $\mathcal{MC}_2$ : for every edge present, the weights of all three types of matrices are non-zero. Formally, the class of matrices  $\mathcal{MC}_2^{>0}$  consists of matrices of the form

$$\mathbf{L}^1 \otimes \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + \mathbf{L}^2 \otimes \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} + \mathbf{L}^{1+2} \otimes \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix},$$

where the non-zero support of all three matrices,  $\mathbf{L}^1$ ,  $\mathbf{L}^2$ , and  $\mathbf{L}^{1+2}$  are the same (but the weights may vary greatly). On the other hand, the matrices that we generate in Section 3.3 can be transformed into such a matrix by adding a small value,  $\delta$  to the weight of all edges where one of the types of edges have non-zero support. We will describe this construction in Section 3.4.1, and bound its condition number in Section 3.4.2.

#### 3.4.1 Construction

In this section, we show the construction from an instance of  $\mathcal{MC}_2$  to an instance of  $\mathcal{MC}_2^{>0}$ . The strategy is to add extra edges with a sufficiently small weight, such that  $\mathbf{L}^1, \mathbf{L}^2, \mathbf{L}^{1+2}$  have identical nonzero structure and the solution of the linear system does not change much.

The reduction from  $\mathcal{MC}_2$  to  $\mathcal{MC}_2^{>0}$ , with pseudo-code in Algorithm 4, simply adds edges with weight  $\delta$  to all the missing edges. The transformation of solutions of the corresponding instance of  $\mathcal{MC}_2^{>0}$  back to a solution of the original  $\mathcal{MC}_2$  instance in Algorithm 5 simply returns the same vector  $\mathbf{x}$ .

---

**Algorithm 4** REDUCE $\mathcal{MC}_2$ TO $\mathcal{MC}_2^{>0}$ 

---

**Input:**  $(B, c^B, \epsilon_1)$  where  $B \in \mathcal{MC}_2$  is an  $m \times n$  matrix,  $c^B \in \mathbb{R}^m$ , and  $0 < \epsilon_1 < 1$ .

**Output:**  $(B^{>0}, c^{B^{>0}}, \epsilon_2)$  where  $B^{>0} \in \mathcal{MC}_2^{>0}$  is an  $m' \times n'$  matrix,  $c^{B^{>0}} \in \mathbb{R}^{m'}$ , and  $0 < \epsilon_2 < 1$ .

- 1:  $\delta \leftarrow \frac{\epsilon_1}{100\kappa(B)\sigma_{\max}(B)\|c^B\|_2}$
- 2:  $\hat{B} \leftarrow \emptyset$  {New equations for  $\mathcal{MC}_2^{>0}$ }
- 3: **for** each pair of vertices  $i, j$  whose blocks are involved in some equation in  $B$  **do**
- 4:   **if**  $B$  does not contain a type 1 equation between  $i$  and  $j$  **then**
- 5:      $\hat{B} \leftarrow \hat{B} \cup \{u_i - u_j = 0\}$ .
- 6:   **end if**
- 7:   **if**  $B$  does not contain a type 2 equation between  $i$  and  $j$  **then**
- 8:      $\hat{B} \leftarrow \hat{B} \cup \{v_i - v_j = 0\}$ .
- 9:   **end if**
- 10:   **if**  $B$  does not contain a type 1 + 2 equation between  $i$  and  $j$  **then**
- 11:      $\hat{B} \leftarrow \hat{B} \cup \{u_i + v_i - (u_j + v_j) = 0\}$ .
- 12:   **end if**
- 13: **end for**
- 14: Let  $\tilde{B}$  be the coefficient matrix of equations in  $\hat{B}$ .
- 15: **return**

$$B^{>0} = \begin{pmatrix} B \\ \delta \tilde{B} \end{pmatrix}, \quad c^{B^{>0}} = \begin{pmatrix} c^B \\ \mathbf{0} \end{pmatrix}, \quad \epsilon_2 = \frac{\epsilon_1}{100}.$$

---

---

**Algorithm 5** MAPSOLN $\mathcal{MC}_2^{>0}$ TO $\mathcal{MC}_2$ 

---

**Input:**  $m \times n$  matrix  $B \in \mathcal{MC}_2$ ,  $m' \times n'$  matrix  $B^{>0} \in \mathcal{MC}_2^{>0}$ , vector  $c^B \in \mathbb{R}^m$  vector  $x \in \mathbb{R}^{n'}$ .

**Output:** Vector  $y \in \mathbb{R}^n$ .

- 1: **if**  $B^\top c^B = \mathbf{0}$  **then**
  - 2:   **return**  $y \leftarrow \mathbf{0}$
  - 3: **else**
  - 4:   **return**  $y \leftarrow x$
  - 5: **end if**
- 

Let  $B^{>0}x = c^{B^{>0}}$  be the linear system returned by a call to REDUCE $\mathcal{MC}_2$ TO $\mathcal{MC}_2^{>0}(B, c^B, \epsilon_1)$ .

As we only add up to two edges per original edge in  $B$ , and do not introduce any new variables, the size of  $B^{>0}$  is immediate from this routine. According to the algorithms,  $c^{B^{>0}}$  is simply

$$c^{B^{>0}} = \begin{pmatrix} c^B \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} c^A \\ \mathbf{0} \end{pmatrix}.$$

Note the two zero vectors in the above equation have different dimensions.

For a given matrix  $\mathbf{B}$ , we will choose the additive term to ensure non-zeros to be:

$$\delta \stackrel{\text{def}}{=} \frac{\epsilon_1}{100\kappa(\mathbf{B})\sigma_{\max}(\mathbf{B})\|\mathbf{c}^{\mathbf{B}}\|_2}. \quad (3.21)$$

Note that we can bound  $\kappa(\mathbf{B}), \sigma_{\max}(\mathbf{B})$  by condition number of  $\mathbf{A}$  (i.e, the linear system instance of  $\mathcal{G}$ ).

**Lemma 3.4.1.**  $\text{nnz}(\mathbf{B}^{>0}) = O(\text{nnz}(\mathbf{B}))$ . *By our setting of  $\delta$  in Equation (3.21), the largest entry of  $\mathbf{B}^{>0}$  does not change, the smallest entry of  $\mathbf{B}^{>0}$  is at least  $\delta$ .*

Note that the addition of all three types of edges means that the null space of  $\mathbf{B}^{>0}$  is now given by the connected components in its graph theoretic structure, and is likely significantly different from the null space of  $\mathbf{B}$ . In order to solve the Linear System Approximation Problem(LSA) problems

$$\begin{aligned} \min_x \|\mathbf{B}\mathbf{x} - \mathbf{c}^{\mathbf{B}}\|_2, \\ \min_x \|\mathbf{B}^{>0}\mathbf{x} - \mathbf{c}^{\mathbf{B}^{>0}}\|_2, \end{aligned}$$

we need to solve the two linear systems:

$$\mathbf{B}^\top \mathbf{B}\mathbf{x} = \mathbf{B}^\top \mathbf{c}^{\mathbf{B}}, \quad (3.22)$$

$$(\mathbf{B}^{>0})^\top \mathbf{B}^{>0}\mathbf{x} = (\mathbf{B}^{>0})^\top \mathbf{c}^{\mathbf{B}^{>0}}. \quad (3.23)$$

First note that the RHS the two equations are the same because:

$$(\mathbf{B}^{>0})^\top \mathbf{c}^{\mathbf{B}^{>0}} = \begin{pmatrix} \mathbf{B}^\top & \delta \tilde{\mathbf{B}}^\top \end{pmatrix} \begin{pmatrix} \mathbf{c}^{\mathbf{B}} \\ \mathbf{0} \end{pmatrix} = \mathbf{B}^\top \mathbf{c}^{\mathbf{B}}.$$

This means that for a sufficiently small choice of  $\delta$ , the differences between the solu-



tions of these two linear systems is small.

**Lemma 3.4.2.** *Let  $\mathbf{B}^{>0} \mathbf{x} = \mathbf{c}^{B^{>0}}$  be the linear system returned by a call to*

$\text{REDUCE}\mathcal{MC}_2\text{TO}\mathcal{MC}_2^{>0}(\mathbf{B}, \mathbf{c}^B, \epsilon_1)$ . *Let  $\mathbf{x}^{B^*} \in \arg \min_{\mathbf{x}} \|\mathbf{B}\mathbf{x} - \mathbf{c}^B\|_2$  and  $\mathbf{x}^{B^*} \perp \text{null}(\mathbf{B})$ .*

*Let  $\mathbf{x}^{B^{>0}*} \in \arg \min_{\mathbf{x}} \|\mathbf{B}^{>0} \mathbf{x} - \mathbf{c}^{B^{>0}}\|_2$ . Then we have:*

$$\left\| \mathbf{x}^{B^*} - \mathbf{x}^{B^{>0}*} \right\|_{\mathbf{B}^\top \mathbf{B}} \leq O\left(\delta(1+\delta)\kappa(\mathbf{B}) \|\mathbf{c}^B\|_2\right).$$

*Proof.* The desired distance can be written as

$$\left\| \mathbf{x}^{B^*} - \mathbf{x}^{B^{>0}*} \right\|_{\mathbf{B}^\top \mathbf{B}} = \left\| \mathbf{B} \left( \mathbf{x}^{B^*} - \mathbf{x}^{B^{>0}*} \right) \right\|_2.$$

Also, the optimality of  $\mathbf{x}^{B^*}$  means that  $\mathbf{B}\mathbf{x}^{B^*} - \mathbf{c}^B$  is perpendicular to anything in the rank space of  $\mathbf{B}$ , in particular,

$$\mathbf{B} \left( \mathbf{x}^{B^*} - \mathbf{x}^{B^{>0}*} \right) \perp \mathbf{B}\mathbf{x}^{B^*} - \mathbf{c}^B,$$

which in turn gives:

$$\left\| \mathbf{B} \left( \mathbf{x}^{B^*} - \mathbf{x}^{B^{>0}*} \right) \right\|_2^2 = \left\| \mathbf{B}\mathbf{x}^{B^{>0}*} - \mathbf{c}^B \right\|_2^2 - \left\| \mathbf{B}\mathbf{x}^{B^*} - \mathbf{c}^B \right\|_2^2. \quad (3.24)$$

We now bound the right hand side. Since  $\mathbf{x}^{B^{>0}*}$  is a minimizer of  $\left\| \mathbf{B}^{>0} \mathbf{x} - \mathbf{c}^{B^{>0}} \right\|_2$ , we have

$$\left\| \mathbf{B}^{>0} \mathbf{x}^{B^{>0}*} - \mathbf{c}^{B^{>0}} \right\|_2^2 \leq \left\| \mathbf{B}^{>0} \mathbf{x}^{B^*} - \mathbf{c}^{B^{>0}} \right\|_2^2,$$

where we can extract out the  $\delta \tilde{\mathbf{B}}$  term in  $\mathbf{B}^{>0}$  separately to get:

$$\left\| \mathbf{B}\mathbf{x}^{B^{>0}*} - \mathbf{c}^B \right\|_2^2 + \delta^2 \left\| \tilde{\mathbf{B}}\mathbf{x}^{B^{>0}*} \right\|_2^2 \leq \left\| \mathbf{B}\mathbf{x}^{B^*} - \mathbf{c}^B \right\|_2^2 + \delta^2 \left\| \tilde{\mathbf{B}}\mathbf{x}^{B^*} \right\|_2^2.$$

Together with Equation (3.24), this then gives

$$\left\| \mathbf{x}^{\mathbf{B}^*} - \mathbf{x}^{\mathbf{B}^{>0*}} \right\|_{\mathbf{B}^\top \mathbf{B}} \leq \delta \left\| \tilde{\mathbf{B}} \mathbf{x}^{\mathbf{B}^*} \right\|_2.$$

It remains to upper bound  $\left\| \tilde{\mathbf{B}} \mathbf{x}^{\mathbf{B}^{>0*}} \right\|_2$ . By the assumption of  $\mathbf{x}^{\mathbf{B}^*} \perp \text{null}(\mathbf{B})$ , we get:

$$\frac{\left\| \tilde{\mathbf{B}} \mathbf{x}^{\mathbf{B}^*} \right\|_2^2}{\left\| \mathbf{B} \mathbf{x}^{\mathbf{B}^*} \right\|_2^2} \leq \frac{\lambda_{\max}(\tilde{\mathbf{B}}^\top \tilde{\mathbf{B}})}{\lambda_{\min}(\mathbf{B}^\top \mathbf{B})}.$$

By a proof similar to Lemma 3.3.8, which upper bounds  $\lambda_{\max}(\mathbf{B}^\top \mathbf{B})$ , we have:

$$\sigma_{\max}(\tilde{\mathbf{B}}) = O(1 + \delta) \sigma_{\max}(\mathbf{B}),$$

which implies

$$\left\| \tilde{\mathbf{B}} \mathbf{x}^{\mathbf{B}^*} \right\|_2 \leq O(1 + \delta) \kappa(\mathbf{B}) \left\| \mathbf{B} \mathbf{x}^{\mathbf{B}^*} \right\|_2.$$

Since  $\mathbf{x}^{\mathbf{B}^*}$  is a minimizer of  $\min_{\mathbf{x}} \left\| \mathbf{B} \mathbf{x} - \mathbf{c}^{\mathbf{B}} \right\|_2$ , we have

$$\left\| \mathbf{B} \mathbf{x}^{\mathbf{B}^*} \right\|_2 = \left\| \Pi_{\mathbf{B}} \mathbf{c}^{\mathbf{B}} \right\|_2 \leq \left\| \mathbf{c}^{\mathbf{B}} \right\|_2.$$

Therefore,

$$\left\| \mathbf{x}^{\mathbf{B}^*} - \mathbf{x}^{\mathbf{B}^{>0*}} \right\|_{\mathbf{B}^\top \mathbf{B}} \leq O(\delta(1 + \delta) \kappa(\mathbf{B})) \left\| \mathbf{c} \right\|_2,$$

which completes the proof. □

We now check that the approximate solutions of the two linear systems in Equation (3.22) and (3.23) are also close to each other.

**Lemma 3.4.3.** *Let  $\mathbf{B}^{>0} \mathbf{x} = \mathbf{c}^{\mathbf{B}^{>0}}$  be the linear system returned by a call to*

*$\text{REDUCE}\mathcal{MC}_2\text{To}\mathcal{MC}_2^{>0}(\mathbf{B}, \mathbf{c}^{\mathbf{B}}, \epsilon_1)$ , and let  $\epsilon_2$  be the error parameter returned by this call.*

Let  $\mathbf{x}$  be a vector such that

$$\left\| \mathbf{B}^{>0} \mathbf{x} - \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{\mathbf{B}^{>0}} \right\|_2 \leq \epsilon_2 \left\| \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{\mathbf{B}^{>0}} \right\|_2.$$

Let  $\mathbf{y}$  be the vector returned by a call to  $\text{MAPSOLNMC}_2^{>0} \text{ToMC}_2(\mathbf{B}, \mathbf{B}^{>0}, \mathbf{c}^{\mathbf{B}}, \mathbf{x})$ . Then

$$\left\| \mathbf{B} \mathbf{y} - \Pi_{\mathbf{B}} \mathbf{c}^{\mathbf{B}} \right\|_2 \leq \epsilon_1 \left\| \Pi_{\mathbf{B}} \mathbf{c}^{\mathbf{B}} \right\|_2.$$

*Proof.* In Algorithm 5, if the condition  $\mathbf{B}^\top \mathbf{c}^{\mathbf{B}} = \mathbf{0}$  is true, then  $\mathbf{y} = \mathbf{0}$ . This implies

$$\left\| \mathbf{B} \mathbf{y} - \Pi_{\mathbf{B}} \mathbf{c}^{\mathbf{B}} \right\|_2 = 0.$$

In the following, we assume that  $\mathbf{B}^\top \mathbf{c}^{\mathbf{B}} \neq \mathbf{0}$ , in which case  $\mathbf{y} = \mathbf{x}$ . We first show that the construction implies  $\left\| \mathbf{B} \mathbf{x} - \Pi_{\mathbf{B}} \mathbf{c}^{\mathbf{B}} \right\|_2 \leq \left\| \mathbf{B}^{>0} \mathbf{x} - \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{\mathbf{B}^{>0}} \right\|_2$ . Once again, let  $\mathbf{x}^{\mathbf{B}^*}$  and  $\mathbf{x}^{\mathbf{B}^{>0}*}$  be the vectors that minimize  $\left\| \mathbf{B} \mathbf{x} - \mathbf{c}^{\mathbf{B}} \right\|_2$  and  $\left\| \mathbf{B}^{>0} \mathbf{x} - \mathbf{c}^{\mathbf{B}^{>0}} \right\|_2$  respectively. This choice gives:

$$\left\| \mathbf{B} \mathbf{x} - \Pi_{\mathbf{B}} \mathbf{c}^{\mathbf{B}} \right\|_2 = \left\| \mathbf{x} - \mathbf{x}^{\mathbf{B}^*} \right\|_{\mathbf{B}^\top \mathbf{B}}.$$

As  $\|\cdot\|_{\mathbf{B}^\top \mathbf{B}}$  is a norm, by triangle inequality we have:

$$\left\| \mathbf{x} - \mathbf{x}^{\mathbf{B}^*} \right\|_{\mathbf{B}^\top \mathbf{B}} \leq \left\| \mathbf{x} - \mathbf{x}^{\mathbf{B}^{>0}*} \right\|_{\mathbf{B}^\top \mathbf{B}} + \left\| \mathbf{x}^{\mathbf{B}^{>0}*} - \mathbf{x}^{\mathbf{B}^*} \right\|_{\mathbf{B}^\top \mathbf{B}}.$$

We will bound these two terms separately.

Since  $\mathbf{B}^\top \mathbf{B} \preceq \mathbf{B}^\top \mathbf{B} + \delta^2 \tilde{\mathbf{B}}^\top \tilde{\mathbf{B}} = (\mathbf{B}^{>0})^\top \mathbf{B}^{>0}$ , the first term is less than its norm in the  $(\mathbf{B}^{>0})^\top \mathbf{B}^{>0}$  norm:

$$\left\| \mathbf{x} - \mathbf{x}^{\mathbf{B}^{>0}*} \right\|_{\mathbf{B}^\top \mathbf{B}} \leq \left\| \mathbf{x} - \mathbf{x}^{\mathbf{B}^{>0}*} \right\|_{(\mathbf{B}^\top \mathbf{B} + \delta^2 \tilde{\mathbf{B}}^\top \tilde{\mathbf{B}})} = \left\| \mathbf{B}^{>0} \mathbf{x} - \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{\mathbf{B}^{>0}} \right\|_2,$$

while the second term is precisely the distances between the two optimums, which we just

bounded in Lemma 3.4.2. Combining these bounds then gives:

$$\|B\mathbf{x} - \Pi_B \mathbf{c}^B\|_2 \leq \left\| B^{>0} \mathbf{x} - \Pi_{B^{>0}} \mathbf{c}^{B^{>0}} \right\|_2 + O\left(\delta(1+\delta)\kappa(\mathbf{B}) \|\mathbf{c}^B\|_2\right). \quad (3.25)$$

As the equations in  $B^{>0}$  is a superset of the ones in  $B$ , we have

$$\left\| B^{>0} \mathbf{x}^{B^{>0}*} - \mathbf{c}^{B^{>0}} \right\|_2 \geq \left\| B \mathbf{x}^{B^{>0}*} - \mathbf{c}^B \right\|_2 \geq \left\| B \mathbf{x}^{B*} - \mathbf{c}^B \right\|_2.$$

Substituting  $\|B \mathbf{x}^{B*} - \mathbf{c}^B\|_2^2 = \|\mathbf{c}^A\|_2^2 - \|\Pi_B \mathbf{c}^B\|_2^2$  and its equivalent in  $B^{>0}$  gives

$$\left\| \Pi_{B^{>0}} \mathbf{c}^{B^{>0}} \right\|_2 \leq \left\| \Pi_B \mathbf{c}^B \right\|_2.$$

Together with the condition of the lemma,

$$\left\| B^{>0} \mathbf{x} - \Pi_{B^{>0}} \mathbf{c}^{B^{>0}} \right\|_2 \leq \epsilon_2 \left\| \Pi_B \mathbf{c}^B \right\|_2.$$

Plugging this into Equation (3.25), we have

$$\|B\mathbf{x} - \Pi_B \mathbf{c}^B\|_2 = O\left(\epsilon_2 \left\| \Pi_B \mathbf{c}^B \right\|_2 + \delta \kappa(\mathbf{B}) \|\mathbf{c}^B\|_2\right).$$

It remains to upper bound  $\|\mathbf{c}^B\|_2$  by a function of  $\|\Pi_B \mathbf{c}^B\|_2$ . By our construction in Algorithm 1 and assumption,  $\|\mathbf{B} \mathbf{c}^B\|_2^2$  is a positive integer. By Lemma 2.2.5,

$$\left\| \Pi_B \mathbf{c}^B \right\|_2 \geq \frac{1}{\sigma_{\max}(\mathbf{B})}.$$

Thus,

$$\|B\mathbf{x} - \Pi_B \mathbf{c}^B\|_2 = O\left(\epsilon_2 + \delta \kappa(\mathbf{B}) \sigma_{\max}(\mathbf{B}) \|\mathbf{c}^B\|_2\right) \left\| \Pi_B \mathbf{c}^B \right\|_2.$$

By our setting of  $\delta$  in Equation (3.21), we have

$$\| \mathbf{B} \mathbf{x} - \Pi_{\mathbf{B}} \mathbf{c}^{\mathbf{B}} \|_2 \leq \epsilon_1 \| \Pi_{\mathbf{B}} \mathbf{c}^{\mathbf{B}} \|_2.$$

This completes the proof. □

### 3.4.2 Bounding Condition Number of the New System in $\mathcal{MC}_2^{>0}$

We now establish bounds on the numerical quantities related to  $\mathbf{B}^{>0}$ . By a proof similar to the upper bound on  $\lambda_{\max}(\mathbf{B}^\top \mathbf{B})$  in Lemma 3.3.8, we have:

$$\lambda_{\max}((\mathbf{B}^{>0})^\top \mathbf{B}^{>0}) = O(\lambda_{\max}(\mathbf{B}^\top \mathbf{B})).$$

As a result, we focus on the lower bound here:

**Lemma 3.4.4.** *The matrix  $\mathbf{B}^{>0}$  from the linear system returned by a call to*

`REDUCE $\mathcal{MC}_2$ TO $\mathcal{MC}_2^{>0}$` ( $\mathbf{B}, \mathbf{c}^{\mathbf{B}}, \epsilon_1, \delta$ ), *where  $\delta$  is set according to Equation 3.21, satisfies*

$$\sigma_{\min}(\mathbf{B}^{>0}) = \Omega\left(\frac{\epsilon_1}{\kappa(\mathbf{B})\sigma_{\max}(\mathbf{B})n_B}\right).$$

*Proof.* Let  $G$  be a unit-edge weight graph whose vertex set and edge set are same as the underlying graph of  $\mathbf{B}^{>0}$ . Let  $\mathbf{L}_G$  be the associated Laplacian matrix of  $G$ , and let  $\mathbf{M} := \mathbf{L}_G \otimes \mathbf{C}$ , where

$$\mathbf{C} = \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}$$

is a symmetric PSD matrix. Note

$$\mathbf{M} = \mathbf{L}_G \otimes \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + \mathbf{L}_G \otimes \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} + \mathbf{L}_G \otimes \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix},$$

which means  $M$  has the same null space as  $(B^{>0})^\top B^{>0}$ .

$$\lambda_{\min}((B^{>0})^\top B^{>0}) = \min_{\mathbf{x} \perp \text{null}(B^{>0})} \frac{\mathbf{x}^\top (B^{>0})^\top B^{>0} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} \geq \delta^2 \min_{\mathbf{x} \perp \text{null}(M)} \frac{\mathbf{x}^\top M \mathbf{x}}{\mathbf{x}^\top \mathbf{x}},$$

the last inequality is due to that  $\delta$  is the minimum edge weight.

Note that the sum of the 3 types of blocks is positive definite and has eigenvalue at least 1: the type 1 and 2 blocks already sum to  $I$ . Formally:

$$\lambda_{\min}(M) = \lambda_{\min}(L_G) \cdot \lambda_{\min}(C) = \lambda_{\min}(L_G).$$

The result then follows from the folklore bound that the minimum non-zero eigenvalue of a unit weighted graph is at least  $\frac{1}{n^2}$ . One way to see this is via Cheeger's inequality (see e.g. [69] applied to each block: this decomposition is equivalent to separating the matrix into its diagonal blocks based on the connected components, and then invoking the fact that the minimum weight of a cut is at least 1, and there are at most  $n$  vertices. Together these tools imply  $\lambda_{\min}(L_G) = \Omega\left(\frac{1}{n_B^2}\right)$ , and in turn:

$$\sigma_{\min}(B^{>0}) \geq \Omega\left(\frac{\delta}{n}\right) = \Omega\left(\frac{\epsilon_1}{\kappa(B)\sigma_{\max}(B)n_B}\right).$$

□

This also implies a bound on the condition number of  $B^{>0}$ :

**Lemma 3.4.5.**  $\kappa(B^{>0}) = O(\epsilon_1^{-1} \sigma_{\max}^2(B) \kappa(B) n_B)$ .

### 3.4.3 Putting it All Together

*Proof of Lemma 3.2.6.* By Lemma 3.4.1, we have

$$\text{nnz}(B^{>0}) = O(s),$$

The largest entry is  $U$ , and the smallest nonzero entry is  $\delta$ . By Equation (3.21) and Lemma 3.1.2 and 3.1.2,

$$\delta = \Omega\left(\frac{\epsilon}{KU^2}\right).$$

By Lemma 3.4.5 and Lemma 3.1.2 and 3.1.2, the condition number is

$$\kappa(\mathbf{B}^{>0}) = O\left(\frac{s^2 U^2 K}{\epsilon}\right).$$

By Lemma 3.4.3, the accuracy error for  $\mathbf{B}^{>0}$  is  $O(\epsilon)$ . □

### 3.5 Rounding and Scaling Weights to Integers

In this section, we show a reduction from the linear system with strict 2-commodity matrix to the linear system with integral strict 2-commodity matrix.

Algorithm 6 presents the pseudo-code for the algorithm  $\text{REDUCE}\mathcal{MC}_2^{>0}\text{To}\mathcal{MC}_{2,\mathbb{Z}}^{>0}$ . Given an instance  $(\mathbf{B}^{>0}, \mathbf{c}^{\mathbf{B}^{>0}}, \epsilon_1)$  where  $\mathbf{B}^{>0} \in \mathcal{MC}_2^{>0}$ , the class of strict 2-commodity matrices, the call  $\text{REDUCE}\mathcal{MC}_2^{>0}\text{To}\mathcal{MC}_{2,\mathbb{Z}}^{>0}(\mathbf{B}^{>0}, \mathbf{c}^{\mathbf{B}^{>0}}, \epsilon_1)$  returns an instance  $(\mathbf{B}^{>0,\mathbb{Z}}, \mathbf{c}^{\mathbf{B}^{>0,\mathbb{Z}}}, \epsilon_2)$  where  $\mathbf{B}^{>0,\mathbb{Z}} \in \mathcal{MC}_{2,\mathbb{Z}}^{>0}$ . Algorithm 7 provides the (trivial) pseudo-code for  $\text{MAPSOLN}\mathcal{MC}_{2,\mathbb{Z}}^{>0}\text{To}\mathcal{MC}_2^{>0}$  which maps a solution of an instance over  $\mathcal{MC}_{2,\mathbb{Z}}^{>0}$  to a solution of an instance over  $\mathcal{MC}_2^{>0}$ .

For simplicity, we analyze an intermediate linear system  $\mathbf{B}^{int} \mathbf{y} = \mathbf{c}^{\mathbf{B}^{>0}}$ , where

$$(\mathbf{B}^{int})_{ij} \stackrel{\text{def}}{=} 2^{-k} \cdot (\mathbf{B}^{>0,\mathbb{Z}})_{ij}.$$

Recall the definition of  $\mathbf{B}^{>0,\mathbb{Z}}$  in Algorithm 6 line 5, we have

$$(\mathbf{B}^{int})_{ij} = 2^{-k} \cdot \lceil (\mathbf{B}^{>0})_{ij} \cdot 2^k \rceil. \tag{3.26}$$

---

**Algorithm 6** REDUCE $\mathcal{MC}_2^{>0}$  TO  $\mathcal{MC}_{2,\mathbb{Z}}^{>0}$ 


---

**Input:**  $(B^{>0}, c^{B^{>0}}, \epsilon_1)$  where  $B^{>0} \in \mathcal{MC}_2^{>0}$  is an  $m \times n$  matrix,  $c^{B^{>0}} \in \mathbb{R}^m$ , and  $0 < \epsilon_1 < 1$

**Output:**  $(B^{>0,\mathbb{Z}}, c^{B^{>0,\mathbb{Z}}}, \epsilon_2)$  where  $B^{>0,\mathbb{Z}} \in \mathcal{MC}_{2,\mathbb{Z}}^{>0}$  is an  $m' \times n'$  matrix,  $c^{B^{>0,\mathbb{Z}}} \in \mathbb{R}^{m'}$ , and  $0 < \epsilon_2 < 1$ .

```

1:  $k \leftarrow \left\lceil \log_2 \frac{\text{nnz}(B^{>0})}{\epsilon_1} \right\rceil + 2$ 
2:  $B^{>0,\mathbb{Z}} \leftarrow B^{>0}$ 
3:  $c^{B^{>0,\mathbb{Z}}} \leftarrow c^{B^{>0}}$ 
4: for each entry  $(B^{>0,\mathbb{Z}})_{ij}$  do
5:    $(B^{>0,\mathbb{Z}})_{ij} \leftarrow \lceil (B^{>0,\mathbb{Z}})_{ij} \cdot 2^k \rceil$ 
6: end for
7: for each entry  $c_i^{B^{>0,\mathbb{Z}}}$  do
8:    $c_i^{B^{>0,\mathbb{Z}}} \leftarrow c_i^{B^{>0,\mathbb{Z}}} \cdot 2^k$ 
9: end for
10: return  $(B^{>0,\mathbb{Z}}, c^{B^{>0,\mathbb{Z}}}, \epsilon_1/3)$ .
```

---



---

**Algorithm 7** MAPSOLN $\mathcal{MC}_{2,\mathbb{Z}}^{>0}$  TO  $\mathcal{MC}_2^{>0}$ 


---

**Input:**  $m \times n$  matrix  $B^{>0} \in \mathcal{MC}_2$ ,  $m' \times n'$  matrix  $B^{>0,\mathbb{Z}} \in \mathcal{MC}_{2,\mathbb{Z}}^{>0}$ , vector  $x \in \mathbb{R}^{n'}$ .

**Output:** Vector  $y \in \mathbb{R}^n$ .

```

1: return  $x$ 
```

---

The linear system  $B^{int}y = c^{B^{>0}}$  is exactly the linear system  $B^{>0,\mathbb{Z}}y = c^{B^{>0,\mathbb{Z}}}$  multiplying a factor  $2^{-k}$  on both sides.

Note the condition number and the eigen-space of a matrix, the optimal solutions of the projection problems are invariant under scaling. Let  $y^* \in \arg \min_y \|B^{>0,\mathbb{Z}}y - c^{B^{>0,\mathbb{Z}}}\|_2 = \arg \min_y \|B^{int}y - c^{B^{>0}}\|_2$ . The following two inequalities:

$$\|y - y^*\|_{(B^{>0,\mathbb{Z}})^\top B^{>0,\mathbb{Z}}} \leq \epsilon \left\| \Pi_{B^{>0,\mathbb{Z}}} c^{B^{>0,\mathbb{Z}}} \right\|_2$$

and

$$\|y - y^*\|_{(B^{int})^\top B^{int}} \leq \epsilon \left\| \Pi_{B^{int}} c^{B^{>0}} \right\|_2$$

are equivalent. Thus, it suffices to analyze the linear system  $B^{int}y = c^{B^{>0}}$ .



Solving the two projection problems

$$\begin{aligned} \min_{\mathbf{x}} \left\| \mathbf{B}^{>0} \mathbf{x} - \mathbf{c}^{\mathbf{B}^{>0}} \right\|_2, \\ \min_{\mathbf{y}} \left\| \mathbf{B}^{int} \mathbf{y} - \mathbf{c}^{\mathbf{B}^{>0}} \right\|_2, \end{aligned}$$

is equivalent to solving the following two linear systems:

$$\begin{aligned} (\mathbf{B}^{>0})^\top \mathbf{B}^{>0} \mathbf{x} &= (\mathbf{B}^{>0})^\top \mathbf{c}^{\mathbf{B}^{>0}}, \\ (\mathbf{B}^{int})^\top \mathbf{B}^{int} \mathbf{y} &= (\mathbf{B}^{int})^\top \mathbf{c}^{\mathbf{B}^{>0}}. \end{aligned} \tag{3.27}$$

Note  $\mathbf{c}^{\mathbf{B}^{>0}} = (\mathbf{c}^{\mathbf{A}}; \mathbf{0})$ , and all entries of the rows of  $\mathbf{B}^{>0}$  corresponding to the original rows of  $\mathbf{A}$  are integers. Thus,

$$(\mathbf{B}^{int})^\top \mathbf{c}^{\mathbf{B}^{>0}} = (\mathbf{B}^{>0})^\top \mathbf{c}^{\mathbf{B}^{>0}}.$$

That is, the 2nd linear system is equivalent to

$$(\mathbf{B}^{int})^\top \mathbf{B}^{int} \mathbf{y} = (\mathbf{B}^{>0})^\top \mathbf{c}^{\mathbf{B}^{>0}}. \tag{3.28}$$

Let

$$\mathbf{M} \stackrel{\text{def}}{=} (\mathbf{B}^{>0})^\top \mathbf{B}^{>0} \text{ and } \widehat{\mathbf{M}} \stackrel{\text{def}}{=} (\mathbf{B}^{int})^\top \mathbf{B}^{int}.$$

We bound the eigenvalues of  $\widehat{\mathbf{M}}$  by the following lemma.

**Lemma 3.5.1.**  $\lambda_{\max}(\widehat{\mathbf{M}}) \leq (1 + 2^{-k+2})\lambda_{\max}(\mathbf{M})$  and  $\lambda_{\min}(\widehat{\mathbf{M}}) \geq \lambda_{\min}(\mathbf{M})$ . Furthermore,  $\kappa(\widehat{\mathbf{M}}) \leq (1 + 2^{-k+2})\kappa(\mathbf{M})$ .

*Proof.* Note  $\mathbf{M}$  can be written as  $\mathbf{B}^\top \mathbf{W} \mathbf{B}$ , where  $\mathbf{B}$  is the incidence-structured block matrix with *unit* nonzero entries and  $\mathbf{W}$  is the diagonal matrix with edge weights. Let  $\mathbf{x}$

be an arbitrary vector, and let  $\mathbf{y} \stackrel{\text{def}}{=} \mathbf{B}\mathbf{x}$ .

$$\mathbf{x}^\top \mathbf{M} \mathbf{x} = \mathbf{y}^\top \mathbf{W} \mathbf{y} = \sum_i \mathbf{W}_{ii} \mathbf{y}_i^2,$$

and

$$\mathbf{x}^\top \widehat{\mathbf{M}} \mathbf{x} = \sum_i \widehat{\mathbf{W}}_i \mathbf{y}_i^2.$$

Let

$$\alpha \stackrel{\text{def}}{=} \min_i \frac{\widehat{\mathbf{W}}_{ii}}{\mathbf{W}_{ii}} \text{ and } \beta \stackrel{\text{def}}{=} \max_i \frac{\widehat{\mathbf{W}}_{ii}}{\mathbf{W}_{ii}}.$$

We have

$$\alpha \sum_i \mathbf{W}_i \mathbf{y}_i^2 \leq \sum_i \widehat{\mathbf{W}}_i \mathbf{y}_i^2 \leq \beta \sum_i \mathbf{W}_i \mathbf{y}_i^2.$$

This implies that

$$\alpha \mathbf{M} \preceq \widehat{\mathbf{M}} \preceq \beta \mathbf{M}. \quad (3.29)$$

Now we bound the values of  $\alpha, \beta$ . According to Equation (3.26),

$$\alpha = \min_i \frac{\widehat{\mathbf{W}}_{ii}}{\mathbf{W}_{ii}} \geq \frac{\left(2^{-k} \cdot \mathbf{W}_{ii}^{1/2} \cdot 2^k\right)^2}{\mathbf{W}_{ii}} = 1. \quad (3.30)$$

Similarly,

$$\beta = \max_i \frac{\widehat{\mathbf{W}}_{ii}}{\mathbf{W}_{ii}} \leq \frac{\left(2^{-k} \left(\mathbf{W}_{ii}^{1/2} \cdot 2^k + 1\right)\right)^2}{\mathbf{W}_{ii}} \leq 1 + 2^{-k+2}. \quad (3.31)$$

The last inequality is due to  $\mathbf{W}_{ii} \geq 1$ . Thus,

$$\lambda_{\max}(\widehat{\mathbf{M}}) \leq (1 + 2^{-k+2}) \lambda_{\max}(\mathbf{M}) \text{ and } \lambda_{\min}(\widehat{\mathbf{M}}) \geq \lambda_{\min}(\mathbf{M}).$$

The bound on the condition number  $\kappa(\widehat{\mathbf{M}}) \leq (1 + 2^{-k+2}) \kappa(\mathbf{M})$  immediately follows the

above two inequalities. □

We show the exact solutions of the two linear systems in Equation (3.27) and (3.28) are close, by the following lemma.

**Lemma 3.5.2.** *Let*

$$\mathbf{x}^* \stackrel{\text{def}}{=} \mathbf{M}^\dagger (\mathbf{B}^{>0})^\top \mathbf{c}^{B^{>0}} \text{ and } \mathbf{y}^* \stackrel{\text{def}}{=} \widehat{\mathbf{M}}^\dagger (\mathbf{B}^{>0})^\top \mathbf{c}^{B^{>0}}.$$

$\mathbf{x}^*$  and  $\mathbf{y}^*$  are exact solutions of linear system (3.27) and (3.28) respectively.

$$\|\mathbf{y}^* - \mathbf{x}^*\|_{\mathbf{M}} \leq 2^{-k+2} \left\| (\mathbf{B}^{>0})^\top \mathbf{c}^{B^{>0}} \right\|_{\mathbf{M}^\dagger}.$$

*Proof.* Note  $\mathbf{M}, \widehat{\mathbf{M}}$  are both symmetric. Expanding the left hand side,

$$\begin{aligned} \|\mathbf{y}^* - \mathbf{x}^*\|_2^2 &= (\mathbf{c}^{B^{>0}})^\top \mathbf{B}^{>0} \left( \widehat{\mathbf{M}}^\dagger - \mathbf{M}^\dagger \right) \mathbf{M} \left( \widehat{\mathbf{M}}^\dagger - \mathbf{M}^\dagger \right) (\mathbf{B}^{>0})^\top \mathbf{c}^{B^{>0}} \\ &= (\mathbf{c}^{B^{>0}})^\top \mathbf{B}^{>0} \mathbf{M}^{\dagger 1/2} \left( \mathbf{M}^{1/2} \widehat{\mathbf{M}}^\dagger \mathbf{M}^{1/2} - \mathbf{I} \right)^2 \mathbf{M}^{\dagger 1/2} (\mathbf{B}^{>0})^\top \mathbf{c}^{B^{>0}} \\ &\leq \left\| \mathbf{M}^{1/2} \widehat{\mathbf{M}}^\dagger \mathbf{M}^{1/2} - \mathbf{I} \right\|_2^2 \left\| \mathbf{M}^{\dagger 1/2} (\mathbf{B}^{>0})^\top \mathbf{c}^{B^{>0}} \right\|_2^2. \end{aligned}$$

The last inequality is by the Courant-Fischer theorem.

By Equation (3.29) and the fact that  $\mathbf{M}, \widehat{\mathbf{M}}$  have same null space,

$$(\beta^{-1} - 1) \mathbf{I} \preccurlyeq \mathbf{M}^{1/2} \widehat{\mathbf{M}}^\dagger \mathbf{M}^{1/2} - \mathbf{I} \preccurlyeq (\alpha^{-1} - 1) \mathbf{I}.$$

By Equation (3.30) and (3.31),

$$\left\| \mathbf{M}^{1/2} \widehat{\mathbf{M}}^\dagger \mathbf{M}^{1/2} - \mathbf{I} \right\|_2 \leq 2^{-k+2}.$$

Therefore,

$$\|\mathbf{y}^* - \mathbf{x}^*\|_M \leq 2^{-k+2} \left\| (\mathbf{B}^{>0})^\top \mathbf{c}^{B^{>0}} \right\|_{M^\dagger}.$$

□

We then show the approximate solutions of the two linear systems in Equation (3.27) and (3.28) are close.

**Lemma 3.5.3.** *Let  $\mathbf{x}$  be a vector such that*

$$\left\| \mathbf{B}^{int} \mathbf{x} - \Pi_{\mathbf{B}^{int}} \mathbf{c}^{B^{>0}} \right\|_2 \leq \epsilon_2 \left\| \Pi_{\mathbf{B}^{int}} \mathbf{c}^{B^{>0}} \right\|_2.$$

Then,

$$\left\| \mathbf{B}^{>0} \mathbf{x} - \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{B^{>0}} \right\|_2 \leq (\epsilon_2 + 2^{-k+2}(1 + \epsilon_2)) \left\| \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{B^{>0}} \right\|_2.$$

*Proof.* Let

$$\mathbf{x}^* \stackrel{\text{def}}{=} \mathbf{M}^\dagger (\mathbf{B}^{>0})^\top \mathbf{c}^{B^{>0}} \text{ and } \mathbf{y}^* \stackrel{\text{def}}{=} \widehat{\mathbf{M}}^\dagger (\mathbf{B}^{>0})^\top \mathbf{c}^{B^{>0}}.$$

Expanding the left hand side,

$$\begin{aligned} \left\| \mathbf{B}^{>0} \mathbf{x} - \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{B^{>0}} \right\|_2 &= \|\mathbf{x} - \mathbf{x}^*\|_M \\ &\leq \|\mathbf{x} - \mathbf{y}^*\|_M + \|\mathbf{y}^* - \mathbf{x}^*\|_M. \end{aligned}$$

The last inequality is due to the triangle inequality. By Equation (3.29), the first term can be upper bounded by

$$\|\mathbf{x} - \mathbf{y}^*\|_M \leq \alpha^{-1/2} \|\mathbf{x} - \mathbf{y}^*\|_{\widehat{\mathbf{M}}}.$$

The second term is upper bounded by Lemma 3.5.2. Thus, we have

$$\left\| \mathbf{B}^{>0} \mathbf{x} - \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{B^{>0}} \right\|_2 \leq \alpha^{-1/2} \epsilon_2 \left\| \Pi_{\mathbf{B}^{int}} \mathbf{c}^{B^{>0}} \right\|_2 + 2^{-k+2} \left\| \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{B^{>0}} \right\|_2.$$

Note that

$$\left\| \Pi_{\mathbf{B}^{int}} \mathbf{c}^{B^{>0}} \right\|_2 = \left\| \mathbf{B}^{>0} \mathbf{c}^{B^{>0}} \right\|_{\widehat{\mathbf{M}}^\dagger} \text{ and } \left\| \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{B^{>0}} \right\|_2 = \left\| \mathbf{B}^{>0} \mathbf{c}^{B^{>0}} \right\|_{\mathbf{M}^\dagger}.$$

Again, by Equation (3.29),

$$\left\| \mathbf{B}^{>0} \mathbf{c}^{B^{>0}} \right\|_{\widehat{\mathbf{M}}^\dagger} \leq \alpha^{-1/2} \left\| \mathbf{B}^{>0} \mathbf{c}^{B^{>0}} \right\|_{\mathbf{M}^\dagger}.$$

By Equation (3.30) and (3.31),

$$\left\| \mathbf{B}^{>0} \mathbf{x} - \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{B^{>0}} \right\|_2 \leq (\epsilon_2 + 2^{-k+2}(1 + \epsilon_2)) \left\| \Pi_{\mathbf{B}^{>0}} \mathbf{c}^{B^{>0}} \right\|_2.$$

This completes the proof. □

*Proof of Lemma 3.2.7.* We set

$$k \stackrel{\text{def}}{=} \left\lceil \log_2 \frac{s}{\epsilon} \right\rceil + 2.$$

After rounding and scaling, the number of nonzero entries does not change. The value of the smallest nonzero entry does not decrease. The value of the largest entry is multiplied by  $2^k = s\epsilon^{-1}$ , which is upper bounded by  $s\epsilon^{-1}U$ .

By Lemma 3.5.1,

$$\kappa(\mathbf{B}^{>0, \mathbb{Z}}) \leq (1 + s^{-1/2}\epsilon^{1/2})\kappa(\mathbf{B}^{>0}).$$

By Lemma 3.5.3, the accuracy is bounded by  $\epsilon/3$ . □

### 3.6 $\mathcal{G}$ Efficiently Reducible to $\mathcal{G}_{\mathbf{L},2}$

In this section, we prove Lemma 3.2.4.

**Definition 3.6.1.** We let  $\mathcal{G}_z$  denote the class of all matrices with integer valued entries s.t. there is at least one non-zero entry in every row and column, and every row has zero row sum.

### 3.6.1 $\mathcal{G} \leq_f \mathcal{G}_z$

In this section, we show the reduction from  $\mathcal{G}$  to  $\mathcal{G}_z$ . Given an instance of linear system  $(\mathbf{A}, \mathbf{c}^A, \epsilon_1)$  with  $\mathbf{A} \in \mathcal{G}$ , the goal is to construct an instance of linear system  $(\mathbf{A}^Z, \mathbf{c}^Z, \epsilon_2)$  with  $\mathbf{A}^Z \in \mathcal{G}_z$  such that, there exists a map between the solutions of the two linear systems.

Algorithm 8 shows the construction of  $(\mathbf{A}^Z, \mathbf{c}^Z, \epsilon_2)$ , and Algorithm 9 shows the transform from a solution of  $(\mathbf{A}^Z, \mathbf{c}^Z, \epsilon_2)$  to a solution of  $(\mathbf{A}, \mathbf{c}^A, \epsilon_1)$ . Recall that  $\mathbf{A}$  has  $n_A$  columns. According to the algorithm,  $\mathbf{A}^Z$  has  $(n_A + 1)$  columns, and  $\mathbf{c}^Z = \mathbf{c}$ .

---

#### Algorithm 8 REDUCE $\mathcal{G}$ TO $\mathcal{G}_z$

---

**Input:**  $(\mathbf{A}, \mathbf{c}^A, \epsilon_1)$  where  $\mathbf{A}_1 \in \mathcal{G}$  is an  $m \times n$  matrix,  $\mathbf{c}^A \in \mathbb{R}^m$ , and  $0 < \epsilon_1 < 1$ .

**Output:**  $(\mathbf{A}^Z, \mathbf{c}^Z, \epsilon_2)$  where  $\mathbf{A}^Z \in \mathcal{G}_z$  is an  $m' \times n'$  matrix,  $\mathbf{c}^Z \in \mathbb{R}^{m'}$ , and  $0 < \epsilon_2 < 1$ .

1: **return**  $((\mathbf{A} \quad -\mathbf{A}\mathbf{1}), \mathbf{c}^A, O(\frac{\epsilon_1}{n}))$ .

---



---

#### Algorithm 9 MAPSOLN $\mathcal{G}_z$ TO $\mathcal{G}$

---

**Input:**  $m \times n$  matrix  $\mathbf{A} \in \mathcal{G}$ ,  $m' \times n'$  matrix  $\mathbf{A}^Z \in \mathcal{G}_z$ , vector  $\mathbf{x}^Z \in \mathbb{R}^{n'}$ .

**Output:** Vector  $\mathbf{x} \in \mathbb{R}^n$ .

1: **return**  $\mathbf{x} \leftarrow \mathbf{x}_{1:n}^Z - \mathbf{x}_{n+1}^Z \mathbf{1}$

---

**Lemma 3.6.2.** *Let  $\mathbf{A}^Z \mathbf{x}^Z = \mathbf{c}^Z$  be the linear system returned by a call to REDUCE  $\mathcal{G}$  TO  $\mathcal{G}_z(\mathbf{A}, \mathbf{c}^A, \epsilon_1)$ . Then,*

$$\text{nnz}(\mathbf{A}^Z) = O(\text{nnz}(\mathbf{A})),$$

*and the largest entry of  $\mathbf{A}^Z$  is at most  $\|\mathbf{A}\|_\infty$ .*

We show the relation between the exact solvers of the two linear systems.

**Claim 3.6.3.** *Let  $\mathbf{x}^{Z*} \in \arg \min_{\mathbf{x}} \|\mathbf{A}^Z \mathbf{x} - \mathbf{c}^Z\|_2$ . Let  $\mathbf{x}^*$  be the output vector of Algorithm 9, that is,*

$$\mathbf{x}^* = \mathbf{x}_{1:n_A}^{Z*} - \mathbf{x}_{n_A+1}^{Z*} \mathbf{1}.$$

Then,

$$\mathbf{x}^* \in \arg \min_x \|\mathbf{A}\mathbf{x} - \mathbf{c}^A\|_2.$$

*Proof.* Note  $\mathbf{1} \in \text{Null}(\mathbf{A}^Z)$ , thus,

$$\mathbf{A}^Z \mathbf{x}^{Z*} = \mathbf{A}^Z (\mathbf{x}^{Z*} - \mathbf{x}_{n_A+1}^{Z*} \mathbf{1}).$$

Expanding it gives

$$\mathbf{A}^Z \mathbf{x}^{Z*} = \begin{pmatrix} \mathbf{A} & -\mathbf{A}\mathbf{1} \end{pmatrix} \begin{pmatrix} \mathbf{x}^* \\ \mathbf{0} \end{pmatrix} = \mathbf{A}\mathbf{x}^*.$$

Note  $\mathbf{c}^Z = \mathbf{c}$ . Therefore,  $\mathbf{x}^* \in \arg \min_x \|\mathbf{A}\mathbf{x} - \mathbf{c}^A\|_2$ . □

We then show the relation between the approximate solvers.

**Lemma 3.6.4.** *Let  $\mathbf{x}^Z$  be a vector such that  $\|\mathbf{A}^Z \mathbf{x}^Z - \Pi_{A^Z} \mathbf{c}^A\|_2 \leq \epsilon_2 \|\Pi_{A^Z} \mathbf{c}^A\|_2$ . Let  $\mathbf{x}$  be the output vector of Algorithm 9, that is,  $\mathbf{x} = \mathbf{x}_{1:n_A}^Z - \mathbf{x}_{n_A+1}^Z \mathbf{1}$ . Then,*

$$\|\mathbf{A}\mathbf{x} - \Pi_A \mathbf{c}^A\|_2 \leq O\left(\frac{\sqrt{n_A} \sigma_{\min}(\mathbf{A})}{\sigma_{\min}(\mathbf{A}^Z)}\right) \epsilon_2 \|\Pi_A \mathbf{c}^A\|_2.$$

*Proof.* Let  $\mathbf{x}^{Z*} \in \arg \min_x \|\mathbf{A}^Z \mathbf{x} - \mathbf{c}^Z\|_2$  whose last entry is 0. Expanding the left hand side norm,

$$\|\mathbf{A}^Z \mathbf{x}^Z - \Pi_{A^Z} \mathbf{c}^A\|_2 = \|\mathbf{A}^Z (\mathbf{x}^Z - \mathbf{x}^{Z*})\|_2.$$

By Claim 3.6.3,  $\mathbf{x}^* \in \arg \min_x \|\mathbf{A}\mathbf{x} - \mathbf{c}^A\|_2$ . Thus,

$$\|\mathbf{A}(\mathbf{x} - \mathbf{x}^*)\|_2 = \|\mathbf{A}^Z(\mathbf{x}^Z - \mathbf{x}^{Z*})\|_2.$$

Note that

$$\mathbf{A}\mathbf{x} = \mathbf{A}^Z \mathbf{x}^Z \text{ and } \mathbf{A}\mathbf{x}^* = \mathbf{A}^Z \mathbf{x}^{Z*}.$$

Together with the Lemma condition, this gives

$$\|\mathbf{A}\mathbf{x} - \Pi_{\mathbf{A}} \mathbf{c}^{\mathbf{A}}\|_2 = \|\mathbf{A}^Z \mathbf{x}^Z - \Pi_{\mathbf{A}^Z} \mathbf{c}^Z\|_2 \leq \epsilon_2 \|\Pi_{\mathbf{A}^Z} \mathbf{c}^{\mathbf{A}}\|_2. \quad (3.32)$$

It remains to upper bound  $\|\Pi_{\mathbf{A}^Z} \mathbf{c}^{\mathbf{A}}\|_2$  by a function of  $\|\Pi_{\mathbf{A}} \mathbf{c}^{\mathbf{A}}\|_2$ . Note

$$\|\Pi_{\mathbf{A}} \mathbf{c}^{\mathbf{A}}\|_2 = \|\mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_{(\mathbf{A}^\top \mathbf{A})^\dagger} \text{ and } \|\Pi_{\mathbf{A}^Z} \mathbf{c}^{\mathbf{A}}\|_2 = \|(\mathbf{A}^Z)^\top \mathbf{c}^{\mathbf{A}}\|_{((\mathbf{A}^Z)^\top \mathbf{A}^Z)^\dagger}. \quad (3.33)$$

It suffices to work on  $\|\mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_{(\mathbf{A}^\top \mathbf{A})^\dagger}$  and  $\|(\mathbf{A}^Z)^\top \mathbf{c}^{\mathbf{A}}\|_{((\mathbf{A}^Z)^\top \mathbf{A}^Z)^\dagger}$ . Since  $\mathbf{A}^Z = \begin{pmatrix} \mathbf{A} & -\mathbf{A}\mathbf{1} \end{pmatrix}$ , we have

$$\|(\mathbf{A}^Z)^\top \mathbf{c}^{\mathbf{A}}\|_2^2 = \|\mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_2^2 + \|\mathbf{1}^\top \mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_2^2.$$

By Cauchy-Schwarz inequality,

$$\|\mathbf{1}^\top \mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_2 \leq \|\mathbf{1}\|_2 \|\mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_2 = \sqrt{n_A} \|\mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_2.$$

Plugging this gives

$$\|(\mathbf{A}^Z)^\top \mathbf{c}^{\mathbf{A}}\|_2^2 \leq (n_A + 1) \|\mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_2^2.$$

By Courant-Fischer theorem,

$$\|(\mathbf{A}^Z)^\top \mathbf{c}^{\mathbf{A}}\|_{((\mathbf{A}^Z)^\top \mathbf{A}^Z)^\dagger}^2 \leq \lambda_{\max} \left( ((\mathbf{A}^Z)^\top \mathbf{A}^Z)^\dagger \right) \|(\mathbf{A}^Z)^\top \mathbf{c}^{\mathbf{A}}\|_2^2 \leq \lambda_{\min}^{-1} ((\mathbf{A}^Z)^\top \mathbf{A}^Z) \|(\mathbf{A}^Z)^\top \mathbf{c}^{\mathbf{A}}\|_2^2$$

and since  $\mathbf{A}^\top \mathbf{c}^{\mathbf{A}}$  is in the eigenspace of  $\mathbf{A}^\top \mathbf{A}$ ,

$$\|\mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_{(\mathbf{A}^\top \mathbf{A})^\dagger}^2 \geq \lambda_{\min} \left( (\mathbf{A}^\top \mathbf{A})^\dagger \right) \|\mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_2^2 = \lambda_{\max}^{-1} (\mathbf{A}^\top \mathbf{A}) \|\mathbf{A}^\top \mathbf{c}^{\mathbf{A}}\|_2^2.$$



Combining the above two inequalities,

$$\|(\mathbf{A}^Z)^\top \mathbf{c}^A\|_{((\mathbf{A}^Z)^\top \mathbf{A}^Z)^\dagger}^2 \leq \frac{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})}{\lambda_{\min}((\mathbf{A}^Z)^\top \mathbf{A}^Z)} (n_A + 1) \|\mathbf{A}^\top \mathbf{c}^A\|_{(\mathbf{A}^\top \mathbf{A})^\dagger}^2.$$

Given Equation (3.32) and (3.33), and condition  $\|\mathbf{A}^Z \mathbf{x}^Z - \mathbf{c}^A\|_2 \leq \epsilon \|\Pi_{\mathbf{A}^Z} \mathbf{c}^A\|_2$ , we have

$$\|\mathbf{A} \mathbf{x} - \Pi_{\mathbf{A}} \mathbf{c}^A\|_2 \leq \epsilon_2 \frac{\sigma_{\max}(\mathbf{A})}{\sigma_{\min}(\mathbf{A}^Z)} \sqrt{n_A + 1} \|\Pi_{\mathbf{A}} \mathbf{c}^A\|_2.$$

This completes the proof.  $\square$

We compute the nonzero condition number of  $\mathbf{A}^Z$ . Note

$$(\mathbf{A}^Z)^\top \mathbf{A}^Z = \begin{pmatrix} \mathbf{M} & -\mathbf{M} \mathbf{1} \\ -\mathbf{1}^\top \mathbf{M} & \mathbf{1}^\top \mathbf{M} \mathbf{1} \end{pmatrix}, \quad (3.34)$$

where  $\mathbf{M} = \mathbf{A}^\top \mathbf{A}$ .

**Claim 3.6.5.**  $\lambda_{\max}((\mathbf{A}^Z)^\top \mathbf{A}^Z) \leq 2n_A \lambda_{\max}(\mathbf{A}^\top \mathbf{A})$ .

*Proof.* Let  $\lambda_1$  be the largest eigenvalue of  $(\mathbf{A}^Z)^\top \mathbf{A}^Z$ , and  $\mathbf{y} = (\tilde{\mathbf{y}}; \alpha)$  be the associated eigenvector of unit length. Let  $\mu_1$  be the largest eigenvalue of  $\mathbf{A}^\top \mathbf{A} = \mathbf{M}$ .

$$\begin{aligned} \lambda_1 &= \mathbf{y}^\top (\mathbf{A}^Z)^\top \mathbf{A}^Z \mathbf{y} = \begin{pmatrix} \tilde{\mathbf{y}}^\top & \alpha \end{pmatrix} \begin{pmatrix} \mathbf{M} & -\mathbf{M} \mathbf{1} \\ -\mathbf{1}^\top \mathbf{M} & \mathbf{1}^\top \mathbf{M} \mathbf{1} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{y}} \\ \alpha \end{pmatrix} \\ &= \tilde{\mathbf{y}}^\top \mathbf{M} \tilde{\mathbf{y}} - 2\alpha \mathbf{1}^\top \mathbf{M} \tilde{\mathbf{y}} + \alpha^2 \mathbf{1}^\top \mathbf{M} \mathbf{1}. \end{aligned}$$

By Courant-Fischer Theorem,

$$\tilde{\mathbf{y}}^\top \mathbf{M} \tilde{\mathbf{y}} \leq \mu_1 \tilde{\mathbf{y}}^\top \tilde{\mathbf{y}}, \text{ and } \mathbf{1}^\top \mathbf{M} \mathbf{1} \leq \mu_1 n_A.$$

By Cauchy-Schwarz inequality,

$$|\mathbf{1}^\top \mathbf{M} \tilde{\mathbf{y}}| \leq \|\mathbf{1}\|_2 \|\mathbf{M} \tilde{\mathbf{y}}\|_2 \leq \sqrt{n_A} \mu_1 \|\tilde{\mathbf{y}}\|_2.$$

Putting all together,

$$\begin{aligned} \lambda_1 &\leq \mu_1 (\|\tilde{\mathbf{y}}\|_2^2 + 2|\alpha| \sqrt{n_A} \|\tilde{\mathbf{y}}\|_2 + \alpha^2 n_A) \\ &= \mu_1 (\|\tilde{\mathbf{y}}\|_2 + |\alpha| \sqrt{n_A})^2. \end{aligned}$$

Since  $\|\tilde{\mathbf{y}}\|_2, |\alpha| \leq 1$ , we have

$$\lambda_1 \leq 2\mu_1 n_A.$$

This completes the proof.  $\square$

Before lower bounding the smallest nonzero singular value of  $\mathbf{A}^Z$ , we characterize the null space of  $\mathbf{A}^Z$  by the following Claim.

**Claim 3.6.6.**

$$\text{null}(\mathbf{A}^Z) = \mathbf{Span} \left( \mathbf{1}, \begin{pmatrix} \mathbf{z} \\ 0 \end{pmatrix} : \mathbf{z} \in \text{null}(\mathbf{A}) \right).$$

*Proof.* Let  $\mathcal{S}$  be the subspace of  $\mathbf{Span} \left( \mathbf{1}, \begin{pmatrix} \mathbf{z} \\ 0 \end{pmatrix} : \mathbf{z} \in \text{null}(\mathbf{A}) \right)$ .

We first show that  $\mathcal{S} \subseteq \text{null}(\mathbf{A}^Z)$ . Clearly,

$$\mathbf{A}^Z \mathbf{1} = \begin{pmatrix} \mathbf{A} & -\mathbf{A} \mathbf{1} \end{pmatrix} \mathbf{1} = \mathbf{0}.$$

For each  $\mathbf{z} \in \text{null}(\mathbf{A})$ , we have

$$\mathbf{A}^Z \begin{pmatrix} \mathbf{z} \\ 0 \end{pmatrix} = \begin{pmatrix} \mathbf{A} & -\mathbf{A} \mathbf{1} \end{pmatrix} \begin{pmatrix} \mathbf{z} \\ 0 \end{pmatrix} = \mathbf{A} \mathbf{z} = \mathbf{0}.$$

Thus,  $\mathcal{S} \subseteq \text{null}(\mathbf{A}^Z)$ .

We then show that  $\mathcal{S} \supseteq \text{null}(\mathbf{A}^Z)$ . For any  $\mathbf{z}' \in \text{null}(\mathbf{A}^Z)$ ,  $\mathbf{z}' - \mathbf{z}'_{n_A+1} \mathbf{1} \in \text{null}(\mathbf{A}^Z)$ . Let  $\mathbf{z} \in \mathbb{R}^n$  such that the  $i$ th entry of  $\mathbf{z}$  is  $\mathbf{z}'_i - \mathbf{z}'_{n_A+1}$ . Thus,

$$\mathbf{0} = \mathbf{A}^Z (\mathbf{z}' - \mathbf{z}'_{n_A+1} \mathbf{1}) = \begin{pmatrix} \mathbf{A} & -\mathbf{A}\mathbf{1} \end{pmatrix} \begin{pmatrix} \mathbf{z} \\ 0 \end{pmatrix} = \mathbf{A}\mathbf{z}.$$

That is, any vector in  $\mathcal{S}$  can be written as a linear combination of  $\mathbf{1}$  and  $\begin{pmatrix} \mathbf{z} \\ 0 \end{pmatrix}$  where  $\mathbf{z} \in \text{null}(\mathbf{A})$ . Thus,  $\mathcal{S} \supseteq \text{null}(\mathbf{A}^Z)$ .

Therefore,  $\text{null}(\mathbf{A}^Z) = \mathcal{S}$ . □

By the above claim, we know that  $\mathbf{A}$  and  $\mathbf{A}^Z$  have same rank. We bound the smallest nonzero singular value of  $\mathbf{A}^Z$  by the following Claim.

**Claim 3.6.7.**  $\lambda_{\min}((\mathbf{A}^Z)^\top \mathbf{A}^Z) \geq \lambda_{\min}(\mathbf{A}^\top \mathbf{A}) / (n_A + 1)$ .

*Proof.* Let  $\lambda_k$  be the smallest nonzero eigenvalue of  $(\mathbf{A}^Z)^\top \mathbf{A}^Z$ , and  $\mathbf{y} = (\tilde{\mathbf{y}}; \alpha)$  be the associated eigenvector of unit length. Let  $\mu_k$  be the smallest nonzero eigenvalue of  $\mathbf{A}^\top \mathbf{A} = \mathbf{M}$ , and  $\mathbf{x}$  be the associated eigenvector of unit length.

$$\lambda_k = \tilde{\mathbf{y}}^\top \mathbf{M} \tilde{\mathbf{y}} - 2\alpha \mathbf{1}^\top \mathbf{M} \tilde{\mathbf{y}} + \alpha^2 \mathbf{1}^\top \mathbf{M} \mathbf{1}.$$

Since  $\mathbf{M}$  is symmetric PSD, it can be written as  $\mathbf{M} = \mathbf{M}^{1/2} \mathbf{M}^{1/2}$ .

$$\lambda_k = \left\| \mathbf{M}^{1/2} (\tilde{\mathbf{y}} - \alpha \mathbf{1}) \right\|_2^2.$$

We decompose  $\tilde{\mathbf{y}} - \alpha \mathbf{1} := \mathbf{z}_1 + \mathbf{z}_2$ , where  $\mathbf{z}_1 \in \text{null}(\mathbf{M})$  and  $\mathbf{z}_2 \perp \text{null}(\mathbf{M})$ . Then,  $\lambda_k = \left\| \mathbf{M}^{1/2} \mathbf{z}_2 \right\|_2^2$ . By Courant-Fischer Theorem

$$\lambda_k \geq \mu_k \|\mathbf{z}_2\|_2^2.$$

To lower bound  $\lambda_k$ , it suffices to lower bound  $\|z_2\|_2$ .

Since  $\mathbf{y} = (\tilde{\mathbf{y}}; \alpha)$ ,

$$\|\tilde{\mathbf{y}} - \alpha \mathbf{1}\|_2^2 = \left\| \mathbf{y} - \alpha \begin{pmatrix} \mathbf{1} \\ 1 \end{pmatrix} \right\|_2^2.$$

Since  $\mathbf{y} \perp \mathbf{1}$ , which is in  $\text{null}(\mathbf{A}^Z)$ , we have

$$\|\tilde{\mathbf{y}} - \alpha \mathbf{1}\|_2^2 = \|\mathbf{y}\|_2^2 + \alpha^2(1 + n_A) = \|\tilde{\mathbf{y}}\|_2^2 + \alpha^2(2 + n_A).$$

Vectors  $\tilde{\mathbf{y}}, -\alpha \mathbf{1}, \tilde{\mathbf{y}} - \alpha \mathbf{1}$  form a triangle. Let  $\theta$  be the angle between  $\tilde{\mathbf{y}}$  and  $\alpha \mathbf{1}$ . See Figure 3.1.

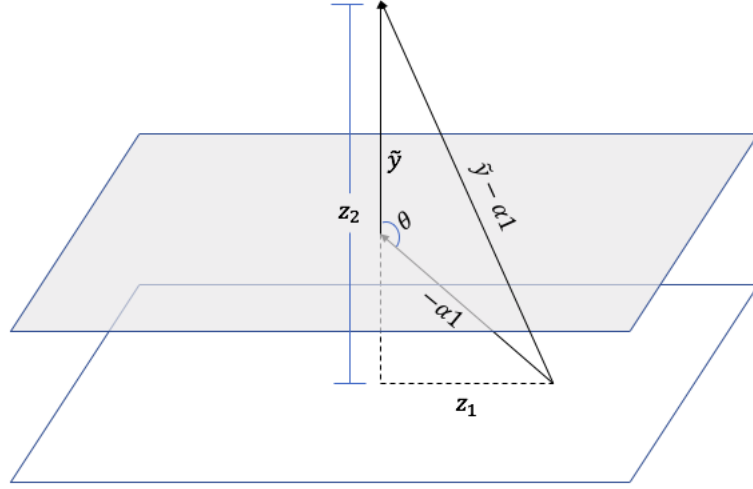


Figure 3.1: The relation between vectors  $\tilde{\mathbf{y}}, -\alpha \mathbf{1}$  and  $\tilde{\mathbf{y}} - \alpha \mathbf{1}$ .

Then,

$$\cos(\theta) = \frac{\|\tilde{\mathbf{y}}\|_2^2 + \|\alpha \mathbf{1}\|_2^2 - \|\tilde{\mathbf{y}} - \alpha \mathbf{1}\|_2^2}{2 \|\tilde{\mathbf{y}}\|_2 \|\alpha \mathbf{1}\|_2} = -\frac{|\alpha|}{\|\tilde{\mathbf{y}}\|_2 \sqrt{n_A}}.$$

Thus,

$$\|z_2\|_2 = \|\tilde{\mathbf{y}}\|_2 + \|\alpha \mathbf{1}\|_2 \cos(\pi - \theta) = \|\tilde{\mathbf{y}}\|_2 + \frac{\alpha^2}{\|\tilde{\mathbf{y}}\|_2} \geq \|\tilde{\mathbf{y}}\|_2.$$

On the other hand,

$$1 = \alpha^2 + \|\tilde{\mathbf{y}}\|_2^2 = \left( \sum_i \tilde{\mathbf{y}}_i \right)^2 + \sum_i \tilde{\mathbf{y}}_i^2 \leq (n_A + 1) \sum_i \tilde{\mathbf{y}}_i^2 = (n_A + 1) \|\tilde{\mathbf{y}}\|_2^2.$$

Therefore,

$$\lambda_k \geq \mu_k \|\tilde{\mathbf{y}}\|_2^2 \geq \frac{\mu_k}{n_A + 1}.$$

This completes the proof. □

The above lemmas give an upper bound on the condition number of  $\mathbf{A}^Z$ .

**Lemma 3.6.8.**  $\kappa(\mathbf{A}^Z) \leq O(n_A^{3/2})\kappa(\mathbf{A})$ .

### 3.6.2 $\mathcal{G}_Z \leq_f \mathcal{G}_{Z,2}$

In this section, we show the reduction from  $\mathcal{G}_Z$  to  $\mathcal{G}_{Z,2}$ . Given an instance of linear system  $(\mathbf{A}^Z, \mathbf{c}^Z, \epsilon_1)$  with  $\mathbf{A}^Z \in \mathcal{G}_Z$ , the goal is to construct an instance of linear system  $(\mathbf{A}^{Z,2}, \mathbf{c}^{Z,2}, \epsilon_2)$  with  $\mathbf{A}^{Z,2} \in \mathcal{G}_{Z,2}$  such that, there is a map between the solutions of these two linear systems.

Algorithm 10 shows the construction of  $(\mathbf{A}^{Z,2}, \mathbf{c}^{Z,2}, \epsilon_2)$ , and Algorithm 11 shows the transform from a solution of  $(\mathbf{A}^{Z,2}, \mathbf{c}^{Z,2}, \epsilon_2)$  to a solution of  $(\mathbf{A}^Z, \mathbf{c}^Z, \epsilon_1)$ . Note  $\mathbf{c}^{Z,2} = (\mathbf{c}^Z; 0)$ .

**Lemma 3.6.9.** *Let  $\mathbf{A}^{Z,2} \mathbf{x}^{Z,2} = \mathbf{c}^{Z,2}$  be the linear system returned by a call to  $\text{REDUCE } \mathcal{G}_Z \text{ TO } \mathcal{G}_{Z,2}(\mathbf{A}^Z, \mathbf{c}^Z, \epsilon_1)$ . Then,*

$$\text{nnz}(\mathbf{A}^{Z,2}) = O(\text{nnz}(\mathbf{A}^Z)).$$

*and the largest entry of  $\mathbf{A}^{Z,2}$  is at most  $\epsilon_1^{-1} \sqrt{m} \sigma_{\max}(\mathbf{A}^Z) \|\mathbf{A}^Z\|_{\infty} \|\mathbf{c}^Z\|_2$ .*

---

**Algorithm 10** REDUCE  $\mathcal{G}_z \text{ TO } \mathcal{G}_{z,2}$ 

---

**Input:**  $(\mathbf{A}^Z, \mathbf{c}^Z, \epsilon_1)$  where  $\mathbf{A}^Z \in \mathcal{G}_z$  is an  $m \times n$  matrix,  $\mathbf{c}^Z \in \mathbb{R}^m$ , and  $0 < \epsilon_1 < 1$ .

**Output:**  $(\mathbf{A}^{Z,2}, \mathbf{c}^{Z,2}, \epsilon_2)$  where  $\mathbf{A}^{Z,2} \in \mathcal{G}_{z,2}$  is an  $m' \times n'$  matrix,  $\mathbf{c}^{Z,2} \in \mathbb{R}^{m'}$ , and  $0 < \epsilon_2 <$

- 1.
  - 1:  $w \leftarrow 3\epsilon_1^{-1}\sqrt{m}\sigma_{\max}(\mathbf{A}^Z) \|\mathbf{A}^Z\|_{\infty} \|\mathbf{c}^Z\|_2$
  - 2:  $k^* \leftarrow \min \{k \in \mathbb{Z} : 2^k \geq \|\mathbf{A}^Z\|_{\infty}\}$
  - 3: Let  $\mathbf{a} \in \mathbb{R}^n$
  - 4: **for**  $i \leftarrow 1$  to  $n$  **do**
  - 5:    $\mathbf{a}_i \leftarrow 2^{k^*} - \|\mathbf{A}_i^Z\|_1 / 2$
  - 6: **end for**
  - 7:  $\mathbf{A}^{Z,2} = \begin{pmatrix} \mathbf{A}^Z & \mathbf{a} & -\mathbf{a} \\ \mathbf{0} & w & -w \end{pmatrix}$
  - 8:  $\mathbf{c}^{Z,2} \leftarrow \begin{pmatrix} \mathbf{c}^Z \\ 0 \end{pmatrix}$
  - 9: **return**  $(\mathbf{A}^{Z,2}, \mathbf{c}^{Z,2}, O(\frac{\epsilon_1}{\sigma_{\max}(\mathbf{A}^Z)\|\mathbf{c}^Z\|_2}))$ .
- 

---

**Algorithm 11** MAPSOLN  $\mathcal{G}_{z,2} \text{ TO } \mathcal{G}_z$ 

---

**Input:**  $m \times n$  matrix  $\mathbf{A}^Z \in \mathcal{G}_z$ ,  $m' \times n'$  matrix  $\mathbf{A}^{Z,2} \in \mathcal{G}_{z,2}$ , vector  $\mathbf{c}^Z \in \mathbb{R}^m$ , vector  $\mathbf{x} \in \mathbb{R}^{n'}$ .

**Output:** Vector  $\mathbf{y} \in \mathbb{R}^n$ .

- 1: **if**  $(\mathbf{A}^Z)^\top \mathbf{c}^Z = \mathbf{0}$  **then**
  - 2:   **return**  $\mathbf{y} \leftarrow \mathbf{0}$
  - 3: **else**
  - 4:   **return**  $\mathbf{y} \leftarrow \mathbf{x}_{1:n}$
  - 5: **end if**
- 

**Claim 3.6.10.**  $\text{null}(\mathbf{A}^{Z,2}) = \text{Span} \left( \begin{pmatrix} \mathbf{0} \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} \mathbf{x} \\ 0 \\ 0 \end{pmatrix} : \mathbf{x} \in \text{null}(\mathbf{A}^Z) \right).$

*Proof.* Let  $\mathbf{y} = (\mathbf{x}; \alpha; \beta) \in \mathbb{R}^{n_A+2}$  such that  $\mathbf{y} \in \text{null}(\mathbf{A}^{Z,2})$ , that is,

$$\mathbf{A}^{Z,2} \mathbf{y} = \begin{pmatrix} \mathbf{A} & \mathbf{a} & -\mathbf{a} \\ \mathbf{0} & w & -w \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} \mathbf{A}^Z \mathbf{x} + (\alpha - \beta) \mathbf{a} \\ (\alpha - \beta) w \end{pmatrix} = \mathbf{0}.$$

Since  $w > 0$ , we have

$$\alpha = \beta.$$

and thus

$$\mathbf{A}^Z \mathbf{x} = \mathbf{0}.$$

Therefore,

$$\text{null}(\mathbf{A}^{Z,2}) = \mathbf{Span} \left( \begin{pmatrix} \mathbf{0} \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} \mathbf{x} \\ 0 \\ 0 \end{pmatrix} : \mathbf{x} \in \text{null}(\mathbf{A}^Z) \right).$$

This completes the proof.  $\square$

We write  $\mathbf{x}^{Z,2}$  as  $(\mathbf{x}; \alpha, \beta)$ . By Claim 3.6.10, the vector  $(\mathbf{0}; 1; 1)$  is in the null space of  $\mathbf{A}^{Z,2}$ . Thus,

$$\mathbf{A}^{Z,2} \mathbf{x}^{Z,2} = \mathbf{A}^{Z,2} \left( \mathbf{x}^{Z,2} - \frac{\alpha + \beta}{2} \begin{pmatrix} \mathbf{0} \\ 1 \\ 1 \end{pmatrix} \right) = \mathbf{A}^{Z,2} \begin{pmatrix} \mathbf{x} \\ \frac{\alpha - \beta}{2} \\ \frac{\beta - \alpha}{2} \end{pmatrix}.$$

Without loss of generality, we assume  $\alpha = \beta$ , that is,  $\mathbf{x}^{Z,2} = (\mathbf{x}; \alpha; -\alpha)$ .

We first show that the exact solutions of the two linear systems are close.

**Lemma 3.6.11.** *Let  $\mathbf{x}^{Z,2*} \in \arg \min_{\mathbf{x}} \|\mathbf{A}^{Z,2} \mathbf{x} - \mathbf{c}^{Z,2}\|_2$ . Write  $\mathbf{x}^{Z,2*} = (\mathbf{s}^*; \alpha; -\alpha)$ . Then,*

$$\mathbf{s}^* \in \arg \min_{\mathbf{x}} \|\mathbf{A}^Z \mathbf{x} - \mathbf{c}^Z\|_P,$$

where  $P = \mathbf{I} - \frac{\mathbf{a}\mathbf{a}^\top}{w^2 + \|\mathbf{a}\|_2^2}$ .

*Proof.* Without loss of generality, we write  $\mathbf{x}^{Z,2} = (\mathbf{x}; \alpha; -\alpha)$ . We expand  $\mathbf{A}^{Z,2}\mathbf{x}^{Z,2} - \mathbf{c}^{Z,2}$ ,

$$\begin{aligned}
& \min_{\mathbf{x}^{Z,2}} \|\mathbf{A}^{Z,2}\mathbf{x}^{Z,2} - \mathbf{c}^{Z,2}\|_2^2 \\
&= \min_{\mathbf{x}, \alpha} \left\| \begin{pmatrix} \mathbf{A}^Z & \mathbf{a} & -\mathbf{a} \\ \mathbf{0} & w & -w \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \alpha \\ -\alpha \end{pmatrix} - \begin{pmatrix} \mathbf{c}^Z \\ 0 \end{pmatrix} \right\|_2^2 \\
&= \min_{\mathbf{x}, \alpha} \|\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z + 2\alpha\mathbf{a}\|_2^2 + 4w^2\alpha^2 \\
&= \min_{\mathbf{x}} \min_{\alpha} 4(w^2 + \|\mathbf{a}\|_2^2) \left( \alpha + \frac{\mathbf{a}^\top(\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z)}{2(w^2 + \|\mathbf{a}\|_2^2)} \right)^2 + \|\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z\|_2^2 - \frac{(\mathbf{a}^\top(\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z))^2}{w^2 + \|\mathbf{a}\|_2^2}.
\end{aligned}$$

Since the minimization over  $\mathbf{x}$  and  $\alpha$  is independent of each other, for any fixed  $\mathbf{x}$ , the above value is minimized when

$$\alpha = -\frac{\mathbf{a}^\top(\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z)}{2(w^2 + \|\mathbf{a}\|_2^2)}.$$

Plugging this value of  $\alpha$  gives

$$\begin{aligned}
\min_{\mathbf{x}^{Z,2}} \|\mathbf{A}^{Z,2}\mathbf{x}^{Z,2} - \mathbf{c}^{Z,2}\|_2^2 &= \min_{\mathbf{x}} \|\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z\|_2^2 - \frac{(\mathbf{a}^\top(\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z))^2}{w^2 + \|\mathbf{a}\|_2^2} \\
&= \min_{\mathbf{x}} (\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z)^\top \left( \mathbf{I} - \frac{\mathbf{a}\mathbf{a}^\top}{w^2 + \|\mathbf{a}\|_2^2} \right) (\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z) \\
&= \min_{\mathbf{x}} \|\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z\|_P^2,
\end{aligned}$$

where  $\mathbf{P} = \mathbf{I} - \frac{\mathbf{a}\mathbf{a}^\top}{w^2 + \|\mathbf{a}\|_2^2}$ . The fact  $\mathbf{x}^{Z,2*} \in \arg \min_{\mathbf{x}} \|\mathbf{A}^{Z,2}\mathbf{x} - \mathbf{c}^{Z,2}\|_2$  implies that

$$\mathbf{s}^* \in \arg \min_{\mathbf{x}} \|\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z\|_P,$$

which completes the proof. □

**Lemma 3.6.12.** Let  $\mathbf{x}^{Z,2*} \in \arg \min_{\mathbf{x}} \|\mathbf{A}^{Z,2}\mathbf{x} - \mathbf{c}^{Z,2}\|_2$  such that  $\mathbf{x}^{Z,2*}$  has the form  $(\mathbf{s}^*; \alpha; -\alpha)$



and  $\|\mathbf{x}^{Z,2*}\|_2$  is minimized. Let  $\mathbf{x}^* \in \arg \min_{\mathbf{x}} \|\mathbf{A}^Z \mathbf{x} - \mathbf{c}^Z\|_2$  such that  $\|\mathbf{x}^*\|_2$  is minimized.

Then,

$$\|\mathbf{A}^Z(\mathbf{s}^* - \mathbf{x}^*)\|_2 \leq \frac{\|\mathbf{a}\|_2}{w} \sigma_{\max}(\mathbf{A}^Z) \|\mathbf{c}^Z\|_2 \|\Pi_{\mathbf{A}^Z} \mathbf{c}^Z\|_2.$$

*Proof.* Let  $\gamma \stackrel{\text{def}}{=} \frac{\|\mathbf{a}\|_2^2}{w^2 + \|\mathbf{a}\|_2^2}$ . By the definition of  $\mathbf{P}$  in Lemma 3.6.11,

$$(1 - \gamma) \mathbf{I} \preceq \mathbf{P} \preceq \mathbf{I}.$$

Thus,

$$(1 - \gamma) \|\mathbf{A}^Z \mathbf{s}^* - \mathbf{c}^Z\|_2^2 \leq \|\mathbf{A}^Z \mathbf{s}^* - \mathbf{c}^Z\|_P^2 \leq \|\mathbf{A}^Z \mathbf{x}^* - \mathbf{c}^Z\|_P^2 \leq \|\mathbf{A}^Z \mathbf{x}^* - \mathbf{c}^Z\|_2^2. \quad (3.35)$$

The second inequality is due to  $\mathbf{s}^* \in \arg \min_{\mathbf{x}} \|\mathbf{A}^Z \mathbf{x} - \mathbf{c}^Z\|_P$ . Note that  $\mathbf{A}^Z \mathbf{x}^* - \mathbf{c}^Z$  is orthogonal to the column space of  $\mathbf{A}^Z$ , in particular, it is orthogonal to  $\mathbf{A}^Z(\mathbf{s}^* - \mathbf{x}^*)$ .

$$\|\mathbf{A}^Z(\mathbf{s}^* - \mathbf{x}^*)\|_2^2 = \|\mathbf{A}^Z \mathbf{s}^* - \mathbf{c}^Z\|_2^2 - \|\mathbf{A}^Z \mathbf{x}^* - \mathbf{c}^Z\|_2^2.$$

By Equation (3.35),

$$\|\mathbf{A}^Z(\mathbf{s}^* - \mathbf{x}^*)\|_2^2 \leq \frac{\gamma}{1 - \gamma} \|\mathbf{A}^Z \mathbf{x}^* - \mathbf{c}^Z\|_2^2 \leq \frac{\gamma}{1 - \gamma} \|\mathbf{c}^Z\|_2^2.$$

If  $\Pi_{\mathbf{A}^Z} \mathbf{c}^Z = \mathbf{0}$ , then  $\mathbf{s}^* = \mathbf{x}^* = \mathbf{0}$  and the claim holds. Otherwise, by Lemma 2.2.5,

$$\|\mathbf{A}^Z(\mathbf{s}^* - \mathbf{x}^*)\|_2 \leq \sqrt{\frac{\gamma}{1 - \gamma}} \sigma_{\max}(\mathbf{A}^Z) \|\mathbf{c}^Z\|_2 \|\Pi_{\mathbf{A}^Z} \mathbf{c}^Z\|_2.$$

Plugging the value of  $\gamma$  completes the proof.  $\square$

We then show the approximate solvers of the two linear systems are close.

**Lemma 3.6.13.** *Let  $\epsilon_2$  be the error parameter returned by a call to REDUCE  $\mathcal{G}_z \text{TO} \mathcal{G}_{z,2}(\mathbf{A}^Z, \mathbf{c}^Z, \epsilon_1)$*

(Algorithm 10). Let  $\mathbf{x}^{Z,2}$  be a vector such that  $\|\mathbf{A}^{Z,2}\mathbf{x}^{Z,2} - \Pi_{\mathbf{A}^{Z,2}}\mathbf{c}^{Z,2}\|_2 \leq \epsilon_2 \|\Pi_{\mathbf{A}^{Z,2}}\mathbf{c}^{Z,2}\|_2$ . Let  $\mathbf{x}$  be the vector returned by Algorithm 11. Then,

$$\|\mathbf{A}^Z\mathbf{x} - \Pi_{\mathbf{A}^Z}\mathbf{c}^Z\|_2 \leq \epsilon_1 \|\Pi_{\mathbf{A}^Z}\mathbf{c}^Z\|_2.$$

*Proof.* If  $\Pi_{\mathbf{A}^Z}\mathbf{c}^Z = \mathbf{0}$ , then  $\mathbf{x} = \mathbf{0}$  and the statement holds. If  $\Pi_{\mathbf{A}^Z}\mathbf{c}^Z \neq \mathbf{0}$ , then  $\mathbf{x} = \mathbf{x}_{1:n_A+1}^{Z,2}$ . Let  $\mathbf{x}^* \in \arg \min_{\mathbf{x}} \|\mathbf{A}^Z\mathbf{x} - \mathbf{c}^Z\|_2$ . We now bound the difference between our solution  $\mathbf{x}$  and  $\mathbf{x}^*$ . Let  $\mathbf{x}^{Z,2*} \in \arg \min_{\mathbf{x}} \|\mathbf{A}^{Z,2}\mathbf{x} - \mathbf{c}^{Z,2}\|_2$  of the form  $(\mathbf{s}^*; \alpha; -\alpha)$ . By the triangle inequality,

$$\|\mathbf{A}^Z(\mathbf{x} - \mathbf{x}^*)\|_2 \leq \|\mathbf{A}^Z(\mathbf{x} - \mathbf{s}^*)\|_2 + \|\mathbf{A}^Z(\mathbf{s}^* - \mathbf{x}^*)\|_2. \quad (3.36)$$

The second term is upper bound by Lemma 3.6.12. It remains to upper bound the first term. Let

$$\delta \stackrel{\text{def}}{=} \epsilon_2 \|\Pi_{\mathbf{A}^{Z,2}}\mathbf{c}^{Z,2}\|_2. \quad (3.37)$$

Without loss of generality, we write  $\mathbf{x}^{Z,2}$  as  $(\mathbf{x}; \beta; -\beta)$ , where  $\mathbf{x}$  is the output of Algo-

rithm 11.

$$\begin{aligned}
\|\mathbf{A}^{Z,2} \mathbf{x}^{Z,2} - \Pi_{\mathbf{A}^{Z,2}} \mathbf{c}^{Z,2}\|_2^2 &= \|\mathbf{A}^{Z,2} (\mathbf{x}^{Z,2} - \mathbf{x}^{Z,2*})\|_2^2 \\
&= \left\| \begin{pmatrix} \mathbf{A}^Z & \mathbf{a} & -\mathbf{a} \\ \mathbf{0} & w & -w \end{pmatrix} \begin{pmatrix} \mathbf{x} - \mathbf{s}^* \\ \beta - \alpha \\ \alpha - \beta \end{pmatrix} \right\|_2^2 \\
&= \left\| \begin{pmatrix} \mathbf{A}^Z (\mathbf{x} - \mathbf{s}^*) + 2(\beta - \alpha) \mathbf{a} \\ 2w(\alpha - \beta) \end{pmatrix} \right\|_2^2 \\
&= \|\mathbf{A} (\mathbf{x} - \mathbf{s}^*) + 2(\beta - \alpha) \mathbf{a}\|_2^2 + 4w^2 (\alpha - \beta)^2.
\end{aligned}$$

Since  $\|\mathbf{A}^{Z,2} \mathbf{x}^{Z,2} - \Pi_{\mathbf{A}^{Z,2}} \mathbf{c}^{Z,2}\|_2 \leq \delta$ , we have

$$4w^2 (\alpha - \beta)^2 \leq \delta^2,$$

that is,

$$(\alpha - \beta)^2 \leq \frac{\delta^2}{4w^2}. \quad (3.38)$$

Similarly, we have

$$\|\mathbf{A}^Z (\mathbf{x} - \mathbf{s}^*) + 2(\beta - \alpha) \mathbf{a}\|_2 \leq \delta.$$

By the triangle inequality,

$$\|\mathbf{A}^Z (\mathbf{x} - \mathbf{s}^*)\|_2 - 2\|(\beta - \alpha) \mathbf{a}\|_2 \leq \delta.$$

Plugging Equation (3.38) and (3.37) into the above inequality, and rearranging it,

$$\|\mathbf{A}^Z (\mathbf{x} - \mathbf{s}^*)\|_2 \leq \left(1 + \frac{\|\mathbf{a}\|_2}{w}\right) \epsilon_2 \|\mathbf{c}^{Z,2}\|_2.$$

By Lemma 2.2.5,

$$\|\mathbf{A}^Z(\mathbf{x} - \mathbf{s}^*)\|_2 \leq \left(1 + \frac{\|\mathbf{a}\|_2}{w}\right) \epsilon_2 \sigma_{\max}(\mathbf{A}^Z) \|\mathbf{c}^{Z,2}\|_2 \|\Pi_{\mathbf{A}^Z} \mathbf{c}^Z\|_2.$$

Together with Lemma 3.6.12 and Equation (3.36), we have

$$\|\mathbf{A}^Z(\mathbf{x} - \mathbf{x}^*)\|_2 \leq \left(\frac{\|\mathbf{a}\|_2}{w} + \left(1 + \frac{\|\mathbf{a}\|_2}{w}\right) \epsilon_2\right) \sigma_{\max}(\mathbf{A}^Z) \|\mathbf{c}^Z\|_2 \|\Pi_{\mathbf{A}^Z} \mathbf{c}^Z\|_2.$$

According to our setting of  $w$  and  $\mathbf{a}$  in line 1 and line 5 of Algorithm 10,

$$w \geq \frac{\|\mathbf{a}\|_2}{\epsilon_2}.$$

This implies that

$$\|\mathbf{A}^Z \mathbf{x} - \Pi_{\mathbf{A}^Z} \mathbf{c}^Z\|_2 \leq \epsilon_1 \|\Pi_{\mathbf{A}^Z} \mathbf{c}^Z\|_2,$$

which completes the proof.  $\square$

We then bound the singular values of  $\mathbf{A}^{Z,2}$ .

**Claim 3.6.14.**  $\lambda_{\max}((\mathbf{A}^{Z,2})^\top \mathbf{A}^{Z,2}) \leq O\left(\epsilon_1^{-2} m \lambda_{\max}((\mathbf{A}^Z)^\top \mathbf{A}^Z) \|\mathbf{A}^Z\|_\infty^2 \|\mathbf{c}^Z\|_2^2\right).$

*Proof.* Let  $\lambda_1$  be the largest eigenvalue of  $(\mathbf{A}^{Z,2})^\top \mathbf{A}^{Z,2}$ , and  $\mathbf{y} = (\tilde{\mathbf{y}}; \alpha; -\alpha)$  be the associated eigenvector of unit length. Let  $\mu_1$  be the largest eigenvalue of  $(\mathbf{A}^Z)^\top \mathbf{A}^Z$ .

$$\begin{aligned} \lambda_1 &= \left\| \begin{pmatrix} \mathbf{A}^Z & \mathbf{a} & -\mathbf{a} \\ \mathbf{0} & w & -w \end{pmatrix} \mathbf{y} \right\|_2^2 \\ &= \left\| \begin{pmatrix} \mathbf{A}^Z & \mathbf{a} & -\mathbf{a} \\ \mathbf{0} & 1 & -1 \end{pmatrix} \mathbf{y} + \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & w-1 & -w+1 \end{pmatrix} \mathbf{y} \right\|_2^2 \\ &= \tilde{\mathbf{y}}^\top (\mathbf{A}^Z)^\top \mathbf{A}^Z \tilde{\mathbf{y}} + 4\alpha \mathbf{a}^\top \mathbf{A}^Z \tilde{\mathbf{y}} + 4\alpha^2 (\mathbf{a}^\top \mathbf{a} + 1) + 8\alpha^2(w-1) + 4\alpha^2(w-1)^2. \end{aligned}$$

By the Courant-Fischer Theorem,

$$\tilde{\mathbf{y}}^\top (\mathbf{A}^Z)^\top \mathbf{A}^Z \tilde{\mathbf{y}} \leq \mu_1 \tilde{\mathbf{y}}^\top \tilde{\mathbf{y}} \leq \mu_1.$$

By Cauchy-Schwarz inequality,

$$|\mathbf{a}^\top \mathbf{A}^Z \tilde{\mathbf{y}}| \leq \|\mathbf{a}\|_2 \|\mathbf{A}^Z \tilde{\mathbf{y}}\|_2 \leq \|\mathbf{A}\|_\infty \sqrt{m} \sqrt{\mu_1}.$$

Thus,

$$\begin{aligned} \lambda_1 &\leq \mu_1 + 4|\alpha| \|\mathbf{A}^Z\|_\infty \sqrt{m} \sqrt{\mu_1} + 4\alpha^2 \left( \|\mathbf{A}^Z\|_\infty^2 m + 1 \right) + 12\alpha^2 w^2 \\ &\leq \left( \sqrt{\mu_1} + 2 \|\mathbf{A}^Z\|_\infty \sqrt{m} \right)^2 + 12\alpha^2 w^2 \\ &\leq 2\mu_1 + 8m \|\mathbf{A}^Z\|_\infty^2 + 12\alpha^2 w^2. \end{aligned}$$

Our setting of  $w$  in line 1 of Algorithm 10 gives

$$\lambda_{\max} \left( (\mathbf{A}^{Z,2})^\top \mathbf{A}^{Z,2} \right) \leq 2\mu_1 + 8m \|\mathbf{A}^Z\|_\infty^2 + 108\epsilon_1^{-2} m \mu_1 \|\mathbf{A}^Z\|_\infty^2 \|\mathbf{c}^Z\|_2^2.$$

This completes the proof. □

**Claim 3.6.15.**  $\lambda_{\min}((\mathbf{A}^{Z,2})^\top \mathbf{A}^{Z,2}) \geq \frac{2\lambda_{\min}((\mathbf{A}^Z)^\top \mathbf{A}^Z)}{2(\|\mathbf{A}^Z\|_\infty^2 m + 1) + \lambda_{\min}((\mathbf{A}^Z)^\top \mathbf{A}^Z)}.$

*Proof.* Let

$$\mathbf{C} = \begin{pmatrix} \mathbf{A}^Z & \mathbf{a} & -\mathbf{a} \\ \mathbf{0} & 1 & -1 \end{pmatrix}.$$

Note

$$\mathbf{A}^{Z,2} = \mathbf{C} + \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & w - 1 & -w + 1 \end{pmatrix}.$$

We can check that

$$(\mathbf{A}^{Z,2})^\top \mathbf{A}^{Z,2} = \mathbf{C}^\top \mathbf{C} + \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & w^2 - 1 & -(w^2 - 1) \\ \mathbf{0} & -(w^2 - 1) & w^2 - 1 \end{pmatrix}.$$

By our setting,  $w^2 - 1 \geq 0$ . The second matrix is a rank-one PSD matrix. Since  $\mathbf{A}^{Z,2}$  and  $\mathbf{C}$  has same null space,

$$\lambda_{\min}((\mathbf{A}^{Z,2})^\top \mathbf{A}^{Z,2}) \geq \lambda_{\min}(\mathbf{C}^\top \mathbf{C}).$$

It suffices to lower bound the smallest nonzero eigenvalue of  $\mathbf{C}^\top \mathbf{C}$ .

Let  $\lambda_k \stackrel{\text{def}}{=} \lambda_{\min}(\mathbf{C}^\top \mathbf{C})$ , and  $\mathbf{y} = (\tilde{\mathbf{y}}; \alpha; -\alpha)$  be the associated eigenvector of unit length. Let  $\mu_k \stackrel{\text{def}}{=} \lambda_{\min}((\mathbf{A}^Z)^\top \mathbf{A}^Z)$ , and  $\mathbf{x}$  be the associated eigenvector of unit length.

$$\begin{aligned} \lambda_k &= \tilde{\mathbf{y}}^\top (\mathbf{A}^Z)^\top \mathbf{A}^Z \tilde{\mathbf{y}} + 4\alpha \mathbf{a}^\top \mathbf{A}^Z \tilde{\mathbf{y}} + 4\alpha^2 (\mathbf{a}^\top \mathbf{a} + 1) \\ &= \|\mathbf{A}^Z \tilde{\mathbf{y}} + 2\alpha \mathbf{a}\|_2^2 + 4\alpha^2 \\ &\geq 4\alpha^2. \end{aligned} \tag{3.39}$$

On the other hand,

$$\begin{aligned} \lambda_k &= \left( 2\alpha \sqrt{\mathbf{a}^\top \mathbf{a} + 1} + \frac{\mathbf{a}^\top \mathbf{A}^Z \tilde{\mathbf{y}}}{\sqrt{\mathbf{a}^\top \mathbf{a} + 1}} \right)^2 + \tilde{\mathbf{y}}^\top (\mathbf{A}^Z)^\top \mathbf{A}^Z \tilde{\mathbf{y}} - \frac{\tilde{\mathbf{y}}^\top (\mathbf{A}^Z)^\top \mathbf{a} \mathbf{a}^\top \mathbf{A}^Z \tilde{\mathbf{y}}}{\mathbf{a}^\top \mathbf{a} + 1} \\ &\geq \tilde{\mathbf{y}}^\top (\mathbf{A}^Z)^\top \left( \mathbf{I} - \frac{\mathbf{a} \mathbf{a}^\top}{\mathbf{a}^\top \mathbf{a} + 1} \right) \mathbf{A}^Z \tilde{\mathbf{y}}. \end{aligned}$$

Take eigen-decomposition of the matrix in the middle,

$$\lambda_k \geq \tilde{\mathbf{y}}^\top (\mathbf{A}^Z)^\top \mathbf{Q} \mathbf{D} \mathbf{Q}^\top \mathbf{A}^Z \tilde{\mathbf{y}},$$

where  $\mathbf{D} := \text{DIAG} \left( 1 - \frac{\mathbf{a}^\top \mathbf{a}}{\mathbf{a}^\top \mathbf{a} + 1}, 1, \dots, 1 \right)$ . By Claim 3.6.10,  $\mathbf{y} \perp \text{null}(\mathbf{A}^{Z,2})$  implies  $\tilde{\mathbf{y}} \perp$

$\text{null}(\mathbf{A}^Z)$ . By the Courant-Fischer Theorem,

$$\lambda_k \geq \frac{1}{\mathbf{a}^\top \mathbf{a} + 1} \|\mathbf{A}^Z \tilde{\mathbf{y}}\|_2^2 \geq \frac{1}{\mathbf{a}^\top \mathbf{a} + 1} \mu_k \|\tilde{\mathbf{y}}\|_2^2.$$

Together with Equation (3.39),

$$\lambda_k \geq \max \left\{ 4\alpha^2, \frac{\mu_k}{\mathbf{a}^\top \mathbf{a} + 1} \|\tilde{\mathbf{y}}\|_2^2 \right\}.$$

Since  $\|\mathbf{y}\|_2^2 = \tilde{\mathbf{y}}_2^2 + 2\alpha^2 = 1$ ,

$$\lambda_k \geq \frac{2\mu_k}{2(\mathbf{a}^\top \mathbf{a} + 1) + \mu_k} \geq \frac{2\mu_k}{2(\|\mathbf{A}\|_\infty^2 m + 1) + \mu_k}.$$

This completes the proof.  $\square$

The above two lemmas give the following bound on the condition number of  $\mathbf{A}^{Z,2}$ .

**Lemma 3.6.16.**  $\kappa(\mathbf{A}^{Z,2}) = O(\epsilon_1^{-1} \sqrt{m} \|\mathbf{A}^Z\|_\infty \|\mathbf{c}^Z\|_2 (\sigma_{\max}(\mathbf{A}^Z) + \kappa(\mathbf{A}^Z) \sqrt{m} \|\mathbf{A}^Z\|_\infty)).$

*Proof of Lemma 3.2.4.* By Lemma 3.6.2 and 3.6.9, we have

$$\text{nnz}(\mathbf{A}^{Z,2}) = O(s).$$

The smallest nonzero entry does not change. The largest entry of  $\mathbf{A}^{Z,2}$  is at most

$$\max\{O(\|\mathbf{A}\|_\infty), w\} = O(\epsilon^{-1} n^2 \sqrt{m} \sigma_{\max}(\mathbf{A}) \|\mathbf{A}\|_\infty \|\mathbf{c}^A\|_2).$$

By Claim 3.1.2 and 3.1.2, the largest entry is upper bounded by

$$O(\epsilon^{-1} s^{9/2} U^3).$$

By Lemma 3.6.8 and 3.6.16, and Claim 3.1.2 and 3.1.2

$$\kappa(\mathbf{A}^{Z,2}) = O(\epsilon^{-1} s^8 U^3 K).$$

By Lemma 3.6.4 and 3.6.13, we have

$$\epsilon_2^{-1} = O(s^{5/2} U^2) \epsilon^{-1}.$$

This completes the proof. □

### 3.7 2D Trusses

In this section, we show that the matrix  $\mathbf{B}$  constructed by  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$  (Algorithm 1), is a 2D Truss Incidence Matrix defined in Definition 2.3.3. It follows that for any function  $f$ ,  $\mathcal{G} \leq_f \mathcal{MC}_2$  implies  $\mathcal{G} \leq_f \mathcal{T}_2$ .

We assume that the algorithm  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$  is called on a matrix  $\mathbf{A}$  with no two identical rows: if there are identical rows, these rows can be collapsed into one row without changing the associated normal equations, by reweighting the resulting row, similar to the technique used in the proof of Claim 3.8.4. Details are left to the reader. A key step is to show that a 2-commodity gadget in the reduction corresponds to a 2D truss subgraph, which we call the 2D-truss gadget.

Without loss of generality, we let  $\mathbf{u}$ -variables correspond to the horizontal axis and  $\mathbf{v}$ -variables to the vertical axis of the 2D plane. According to Definition 2.3.1 and 2.3.3:

1. an equation  $\mathbf{u}_i - \mathbf{u}_j = 0$  in a 2-commodity linear system corresponds to a horizontal edge in the 2D plane;
2. an equation  $\mathbf{v}_i - \mathbf{v}_j = 0$  in a 2-commodity linear system corresponds to a vertical edge in the 2D plane;



3. an equation  $u_i - v_i - (u_j - v_j) = 0$  in a 2-commodity linear system corresponds to a diagonal edge in the 2D plane.

Note that our reduction here heavily relies on the ability to choose arbitrary weights. In particular, the weights on the elements are not related at all with the distances between the corresponding vertices.

Our strategy for picking the coordinates of the vertices of the constructed 2D truss is the following: we first pick the coordinates of the original  $n$  vertices randomly, and then determine the coordinates of the new vertices constructed in the reduction to satisfy all the truss equations.

For the  $n$  original vertices, we pick their  $u$ -coordinates arbitrarily and pick their  $v$ -coordinates randomly. We pick an  $n$ -dimensional random vector  $\mathbf{y}$  uniformly distributed on the  $n$ -dimensional sphere centered at the origin and with radius  $R = \|\mathbf{A}\|_1 n^{10}$ . We then round each entry of  $\mathbf{y}$  to have precision  $\delta = 10^{-10}$ , so that the total number of bits used to store an entry is at most  $O(\log(n \|\mathbf{A}\|_1))$ . Let  $\tilde{\mathbf{y}}$  be the vector after rounding. We assign the  $v$ -coordinate of the  $i$ th vertex to be the  $i$ th entry of  $\tilde{\mathbf{y}}$ .

We then pick the coordinates of the new vertices in the order they are created. Note that each time we replace two vertices in the current equations, say  $s^{j_1}, s^{j_2}$ , whose coordinates have already been determined, we create a 2D truss gadget with 7 new vertices, say  $s^t, s^{t+1}, \dots, s^{t+6}$  (See Algorithm 2  $\mathcal{MC}_2\text{GADGET}$  for the construction.). According to the construction of this gadget, the new vertices  $s^{t+1}, \dots, s^{t+6}$  only appear in this single gadget, whose coordinates do not affect other vertices. Figure 3.2 is the corresponding subgraph which satisfies all the equations in the 2D truss gadget. Note the two triangles  $(s^{t+3}, s^{t+5}, s^{t+6})$  and  $(s^{t+3}, s^{t+4}, s^{t+5})$  need to be isosceles right triangles, which implies  $v_t = (v_{j_1} + v_{j_2})/2$ . Note also that we can assign  $u$ -coordinates to the new vertices which are not between the  $u$ -coordinates of  $s^{j_1}$  and  $s^{j_2}$ . In fact, using an appropriate choice of  $u$ -coordinates and edge weights, we can always place  $s^t, s^{t+1}, \dots, s^{t+6}$  to get the desired equations, provided  $v_{j_1} \neq v_{j_2}$ , which we later will argue holds with high probability using

Lemma 3.7.1. First, however, we state and prove Lemma 3.7.1.

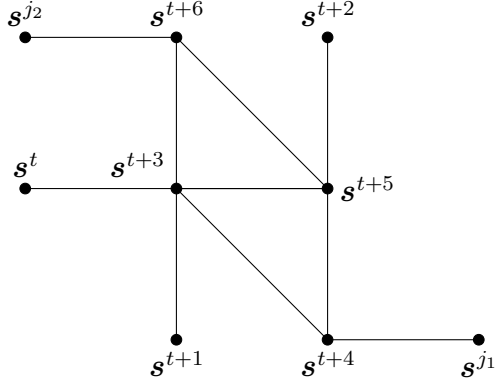


Figure 3.2: Geometric realization of the multicommodity flow gadget generated by as a truss matrix

**Lemma 3.7.1.** *Let  $\mathbf{a} \in \mathbb{R}^n$  be a fixed vector such that  $-2 \leq \mathbf{a}_i \leq 2, \forall i \in [n]$  and  $\mathbf{a}^\top \mathbf{1} = 0$ , and  $\mathbf{a} \neq \mathbf{0}$ . Let  $\tilde{\mathbf{y}}$  be a vector picked as above. Then,*

$$\Pr(\mathbf{a}^\top \tilde{\mathbf{y}} = 0) \leq \frac{2\delta n^2}{\|\mathbf{a}\|_2 R}.$$

*Proof.* Let  $\Delta \stackrel{\text{def}}{=} \tilde{\mathbf{y}} - \mathbf{y}$ . Clearly,  $-\delta \leq \Delta_i \leq \delta, \forall i \in [n]$ .

$$\mathbf{a}^\top \tilde{\mathbf{y}} = \mathbf{a}^\top (\mathbf{y} + \Delta) \leq \mathbf{a}^\top \mathbf{y} + \sum_{i \in [n]} |\mathbf{a}_i| |\Delta_i| \leq \mathbf{a}^\top \mathbf{y} + 2\delta n.$$

Similarly,  $\mathbf{a}^\top \tilde{\mathbf{y}} \geq \mathbf{a}^\top \mathbf{y} - 2\delta n$ . Thus,

$$\Pr(\mathbf{a}^\top \tilde{\mathbf{y}} = 0) \leq \Pr(|\mathbf{a}^\top \mathbf{y}| \leq 2\delta n).$$

Since the distribution of  $\mathbf{y}$  is rotation invariant, we assume without loss of generality  $\mathbf{a} =$

$(\|\mathbf{a}\|_2, 0, \dots, 0)$ . Let  $A$  be the area of the  $n$ -dimensional sphere.

$$A = \frac{2\pi^{(n+1)/2}}{\Gamma\left(\frac{n+1}{2}\right)} \cdot R^n.$$

Let  $A_\delta$  be the area of  $\{\mathbf{y} : \|\mathbf{y}\|_2 = R, |\mathbf{y}_1| \leq 2\delta n / \|\mathbf{a}\|_2\}$ . Then,

$$A_\delta \leq \frac{2\pi^{n/2}}{\Gamma\left(\frac{n}{2}\right)} \cdot R^{n-1} \cdot \frac{2\delta n}{\|\mathbf{a}\|_2}.$$

Thus, (assume  $n$  is even, the case of odd  $n$  can be checked similarly)

$$\begin{aligned} \Pr(|\mathbf{a}^\top \mathbf{y}| \leq \delta) &= \frac{A_\delta}{A} \\ &\leq \frac{2\delta n}{\|\mathbf{a}\|_2 \sqrt{\pi} R} \cdot \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)} \\ &= \frac{2\delta n}{\|\mathbf{a}\|_2 R} \cdot \frac{n! \sqrt{\pi}}{n!!(n-2)!!} \\ &\leq \frac{2\delta n^2}{\|\mathbf{a}\|_2 R}. \end{aligned}$$

This completes the proof. □

By our construction of the truss, for each vertex, its  $\mathbf{v}$ -coordinate can be written as a fixed convex combination of  $\tilde{\mathbf{y}}$ , say  $\mathbf{c}^\top \tilde{\mathbf{y}}$  in which  $\mathbf{c}^\top \mathbf{1} = 1$  and  $c_i \geq 0, \forall i \in [n]$ . Note that when algorithm  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$  is applied to a matrix  $\mathbf{A}$ , it processes the  $j$ th row in  $k$  iterations where by Lemma 3.3.1 we have  $k \leq \log \|\mathbf{A}_j\|_1$ . Let  $p \stackrel{\text{def}}{=} \log \|\mathbf{A}_j\|_1$ . Given how the convex combination specified by  $\mathbf{c}$  is formed, it follows that  $2^p \mathbf{c}$  is an integer vector.

Next we argue that when two variables are chosen for pairing by the  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$  algorithm as it processes some row  $\mathbf{A}_j$ , these two variables will have their  $\mathbf{v}$ -coordinates represented as convex combinations  $\mathbf{c}^\top \tilde{\mathbf{y}}$  and  $\mathbf{d}^\top \tilde{\mathbf{y}}$  where crucially  $\mathbf{c} \neq \mathbf{d}$ . This ensures that  $2^p(\mathbf{c} - \mathbf{d})$  is a non-zero integer vector, and hence  $\|\mathbf{c} - \mathbf{d}\|_2 \geq 2^{-p} = 1 / \|\mathbf{A}_j\|_1$ . This follows from stronger claim stated below, which we prove later.

**Claim 3.7.2.** Suppose algorithm  $\text{REDUCE } \mathcal{G}_{z,2}\text{TO}\mathcal{MC}_2$  is processing some row  $\mathbf{A}_j$ . Let  $V^l$  be the set of variables with non-zero coefficients in the main equation  $\mathcal{A}_j$  at the  $l$ th iteration of the while-loop in Line 20 of  $\text{REDUCE } \mathcal{G}_{z,2}\text{TO}\mathcal{MC}_2$ . Let  $S^l$  be the set of associated vertices. Consider two arbitrary vertices  $\mathbf{s}, \mathbf{t} \in S^l$ , and let  $\mathbf{c}, \mathbf{d} \in \mathbb{R}^n$  be the non-negative vectors with  $\mathbf{c}^\top \mathbf{1} = \mathbf{d}^\top \mathbf{1} = 1$  s.t. the  $\mathbf{v}$ -coordinates of  $\mathbf{s}$  and  $\mathbf{t}$  represented as convex combinations are  $\mathbf{c}^\top \tilde{\mathbf{y}}$  and  $\mathbf{d}^\top \tilde{\mathbf{y}}$  respectively. For every  $i \in [n]$ , we view entries  $\mathbf{c}_i$  and  $\mathbf{d}_i$  as fixed point binary numbers:  $\mathbf{c}_i = \alpha_0^i \alpha_1^i \alpha_2^i \dots \alpha_p^i$  and  $\mathbf{d}_i = \beta_0^i \beta_1^i \beta_2^i \dots \beta_p^i$ , or equivalently  $\mathbf{c}_i = \sum_{k=0}^p \alpha_k^i 2^{-k}$  and  $\mathbf{d}_i = \sum_{k=0}^p \beta_k^i 2^{-k}$ . Then there is no index  $k$  s.t.  $1 = \alpha_k^i = \beta_k^i$ , i.e. there is no index where the  $k$ th bit is 1 in both strings.

These two vertices have same  $\mathbf{v}$ -coordinate if and only if  $(\mathbf{c} - \mathbf{d})^\top \tilde{\mathbf{y}} = 0$ . Let  $\mathbf{a} \stackrel{\text{def}}{=} \mathbf{c} - \mathbf{d}$ . Then,  $-2 \leq \mathbf{a}_i \leq 2, \forall i \in [n]$ ,  $\mathbf{a}^\top \mathbf{1} = 0$ , and

$$\|\mathbf{a}\|_2 \geq 1 / \|\mathbf{A}_j\|_1.$$

By Lemma 3.7.1,

$$\Pr(\mathbf{c}^\top \tilde{\mathbf{y}} = \mathbf{d}^\top \tilde{\mathbf{y}}) \leq \frac{2\delta n^2 \|\mathbf{A}_j\|_1}{R}.$$

By Lemma 3.3.1, the total number of the vertices in the truss is at most

$$O(n^2 \log n).$$

By a union bound, the probability that there exist two different vertices with same  $\mathbf{v}$ -coordinate is at most

$$\frac{2\delta n^2 \|\mathbf{A}\|_1}{R} \cdot O(n^2 \log n)^2 = O\left(\frac{\log^2 n}{n^4}\right).$$

*Proof of Lemma 3.2.8.* Since the linear system for 2D trusses is the same as the linear system for 2-commodity, all complexity parameters of these two linear systems are the same. □

*Proof of Claim 3.7.2.* We prove the claim by induction. Our induction hypothesis is simply that the claim holds in round  $l$ .

Note that by Lemma 3.3.1 every time a new vertex is created (during the  $l$ th iteration of while-loop in Line 20), it is always paired in the following iteration (iteration  $l + 1$ ) and then disappears (i.e. has zero coefficient) in the main equation in all following iterations.

Observe that the convex combination vector  $\mathbf{a}$  for each vertex that corresponds to an original variable  $i$  is has  $\mathbf{a}_i = 1$  and  $\mathbf{a}_h = 0$  for all  $h \neq i$ . This proves the induction hypothesis for the base case of the variables in the main equation  $\mathcal{A}_j$  before the first iteration of the while-loop (i.e.  $l = 0$ ).

Suppose  $\mathbf{c}, \mathbf{d}$  are the convex combination vectors for two variables that exist in some round  $l$ . Assume the induction hypothesis holds for round  $l - 1$ . Write the binary strings for the  $i$ th of both vectors as  $\mathbf{c}_i = \alpha_0^i \alpha_1^i \alpha_2^i \dots \alpha_p^i$  and  $\mathbf{d}_i = \beta_0^i \beta_1^i \beta_2^i \dots \beta_p^i$ , or equivalently  $\mathbf{c}_i = \sum_{k=0}^p \alpha_k^i 2^{-k}$  and  $\mathbf{d}_i = \sum_{k=0}^p \beta_k^i 2^{-k}$ . Suppose for the sake of contradiction that there exists some  $k$  s.t.  $\alpha_k^i = \beta_k^i = 1$ . Trivially, it cannot be the case that the such a collision occurs if either variable is an original variable. So both variables must be new variables. The bit string for entry  $\mathbf{c}_i$  is created by averaging two bit strings of variables from the main equation in round  $l - 1$ , say  $\eta_0^i \eta_1^i \eta_2^i \dots \eta_p^i$  and  $\gamma_0^i \gamma_1^i \gamma_2^i \dots \gamma_p^i$ . Similarly, bit string for entry  $\mathbf{d}_i$  is created by averaging two bit strings of variables from the main equation in round  $l - 1$ , say  $\theta_0^i \theta_1^i \theta_2^i \dots \theta_p^i$  and  $\sigma_0^i \sigma_1^i \sigma_2^i \dots \sigma_p^i$ . Note that each variable can only be paired once in each iteration, so the four bit strings must come from distinct variables in round  $l$ .

$\alpha_k^i = 1$  requires that exactly one of the following conditions is true:

1.  $\gamma_{k-1}^i = 1$
2.  $\eta_{k-1}^i = 1$
3. A “carry” occurred when adding strings  $\gamma_k^i \gamma_{k+1}^i \dots \gamma_p^i$  and  $\eta_k^i \eta_{k+1}^i \dots \eta_p^i$ .

But, Case 3 immediately leads to a contradiction, as a carry can only occur when there exists some bit position  $g$  s.t.  $\gamma_g^i = \eta_g^i$ . But this is false by the induction hypothesis. Thus

we must be in either Case 1 or Case 2. By similar logic, we can conclude from  $\beta_k^i = 1$  that exactly one of the following must be true: either  $\theta_{k-1}^i = 1$  or  $\sigma_{k-1}^i = 1$ . All together, we have concluded that exactly two of the four bits  $\eta_{k-1}^i, \gamma_{k-1}^i, \theta_{k-1}^i$ , and  $\sigma_{k-1}^i$  must be set to 1. This contradicts the induction hypothesis. Having established a contradiction whenever the induction hypothesis fails at step  $l$ , we have shown that it holds at this step.  $\square$

### 3.8 Connections with Interior Point Methods

In this section, we discuss how applications in scientific computing and combinatorial optimization produce the linear systems that we show are hard to solve. We first give a brief overview of interior point methods, with focus on how they generate linear systems in Section 3.8.1. Then we formalize the matrices that interior point methods produce when run on 2-commodity flow matrices in Section 3.8.2 and isotropic total variation minimization in Section 3.8.3.

#### 3.8.1 Brief Overview of Interior Point Methods

Interior point methods [70, 71, 72, 73, 24, 4, 74] can be viewed as ways of solving convex optimization problems via a sequence of linear systems. For simplicity, we choose the log-barrier based interpretation from Chapter 11 of the book by Boyd and Vandenberghe as our starting point. The main idea is to represent a convex optimization problem with constraints as a sequence of linear programs with terms called barrier functions added to the objective. Specifically, turning the problem:

$$\begin{aligned} \min_{\mathbf{y}} \quad & f(\mathbf{y}) \\ \text{subject to:} \quad & \mathbf{M}\mathbf{y} = \mathbf{b} \\ & \mathbf{y} \geq 0 \end{aligned}$$

into the equality-constrained optimization problem:

$$\begin{aligned} \min_{\mathbf{y}} \quad & t \cdot f(\mathbf{y}) - \sum_i \log(\mathbf{y}_i) \\ \text{subject to:} \quad & \mathbf{M}\mathbf{y} = \mathbf{b} \end{aligned}$$

for a parameter  $t$  that is gradually increased (by factors of about  $1 + n^{-1/2}$  throughout the course of the algorithm). The solution of the above optimization problem converges to an optimal solution of the original linear programming, as  $t$  goes to infinity. Between these increase steps, the algorithm performs Newton steps on this log barrier objective, which when combined with the equality constraint  $\mathbf{M}\mathbf{y} = \mathbf{b}$  requires solving the problem:

$$\begin{aligned} \min_{\Delta \mathbf{y}} \quad & -\mathbf{g}(\mathbf{y})^\top \Delta \mathbf{y} \\ \text{subject to:} \quad & \mathbf{M}\Delta \mathbf{y} = \mathbf{0} \\ & \|\Delta \mathbf{y}\|_{H(\mathbf{y})} \leq 0.1 \end{aligned}$$

which is to maximize the projection along the gradient subject to the second order term being at most 0.1 and staying in the null space. This can in turn be interpreted as a least squares problem, and solving the linear system:

$$\mathbf{M}\mathbf{H}(\mathbf{y})^{-1}\mathbf{M}^\top \mathbf{x} = -t\mathbf{M}\mathbf{H}(\mathbf{y})^{-1}\mathbf{g}(\mathbf{y}). \quad (3.40)$$

### 3.8.2 2-Commodity Flow

We now show that solving 2-commodity flow problems using interior point methods as described in Subsection 3.8.1 can lead to any system in the class  $\mathcal{MC}_2^{>0}$ . There are also many variants of the multicommodity flow problem [75], and we work with the minimum cost version due to it being the most general.

**Definition 3.8.1** (Min-cost 2-commodity flow problem). *Given a directed graph  $G =$*

$(V, E)$  with  $n$  vertices and  $m$  edges, a positive edge-capacity vector  $\mathbf{z} \in \mathbb{R}^m$ , a positive edge-cost vector  $\mathbf{c} \in \mathbb{R}^{2m}$ , and two vertex-demand vectors  $\mathbf{d}^1, \mathbf{d}^2 \in \mathbb{R}^n$ . The goal is to compute two flows  $\mathbf{y}^1, \mathbf{y}^2$  such that,

1.  $\mathbf{y}^1$  satisfies the demand  $\mathbf{d}^1$ , and  $\mathbf{y}^2$  satisfies the demand  $\mathbf{d}^2$ ,
2. for each edge  $e \in E$ , the sum of the two flows on  $e$  is no larger than the edge capacity  $z_e$ , and
3. the total cost of the two flows is minimized.

To formulate this as a linear program, we let  $\mathbf{y} = (\mathbf{y}^1; \mathbf{y}^2)$  be the two flows, and  $\mathbf{N}$  be the edge-vertex incidence matrix of graph  $G$ .

$$\begin{aligned}
\min \quad & \mathbf{c}^\top \mathbf{y} \\
s.t. \quad & \mathbf{N}^\top \mathbf{y}^1 = \mathbf{d}^1 \\
& \mathbf{N}^\top \mathbf{y}^2 = \mathbf{d}^2 \\
& \mathbf{y}^1, \mathbf{y}^2, \mathbf{z} - \mathbf{y}^1 - \mathbf{y}^2 \geq \mathbf{0}
\end{aligned}$$

Write this linear programming as the following minimization problem with a logarithmic barrier function,

$$\begin{aligned}
\min \quad & \mathbf{c}^\top \mathbf{y} - \frac{1}{t} \sum_{i \in [m]} \log \mathbf{y}_i^1 + \log \mathbf{y}_i^2 + \log(\mathbf{z}_i - \mathbf{y}_i^1 - \mathbf{y}_i^2) \\
s.t. \quad & \mathbf{N}^\top \mathbf{y}^1 = \mathbf{d}^1 \\
& \mathbf{N}^\top \mathbf{y}^2 = \mathbf{d}^2
\end{aligned}$$

where  $t > 0$  is a parameter.

**Definition 3.8.2.** We say a linear system  $\mathbf{B}^\top \mathbf{B} \mathbf{x} = \mathbf{B}^\top \mathbf{c}$  is a Minimum Cost 2-commodity Flow IPM Linear System if it can be obtained from as an instance of the Newton-Step Linear System of Equation (3.40) for some Min-cost 2-commodity flow problem.



The class of Minimum Cost 2-commodity Flow IPM Linear Systems is appears more restrictive than  $\mathcal{MC}_2^{>0}$ , but as we will see, it is essentially equivalent, and still sufficiently expressive that any linear system can be reduced to it. The main result that we will sketch in this section is:

**Lemma 3.8.3.** *For any linear system  $\mathbf{Ax} = \mathbf{c}$  (with error parameter  $\epsilon$ ) with polynomially bounded sparse parameter complexity, there exists an efficient reduction to a Minimum Cost 2-commodity Flow IPM Linear System.*

Our reduction from a general linear system  $(\mathbf{A}, \mathbf{c}, \epsilon)$ , where  $\mathbf{A} \in \mathcal{G}$ , to a Minimum Cost 2-commodity Flow IPM Linear System is the same as our reduction to  $\mathcal{MC}_2^{>0}$ , except that before our chain of reductions, we first multiply the linear system by a diagonal matrix  $\mathbf{S}$  with diagonal entries that are  $\pm 1$ . We will later specify how to choose these signs. This gives us the linear system  $\mathbf{SAx} = \mathbf{Sc}$  with the same error parameter as before. It is easy to verify that this system has the same sparse parameter complexity as the original linear system, and an approximate solution to this system is an approximate solution to the original system with the same  $\epsilon$ .

We now apply our usual chain of reductions to get a linear system in over a matrix in  $\mathcal{MC}_2^{>0}$ . We write this system as

$$(\mathbf{B}^{>0})^\top \mathbf{B}^{>0} \mathbf{x} = (\mathbf{B}^{>0})^\top \mathbf{c}^{\mathbf{B}^{>0}}.$$

The remainder of this Section is dedicated to showing that this linear system is a Minimum Cost 2-commodity Flow IPM Linear System, thus proving Lemma 3.8.3.

We can pull out the edge weights from  $\mathbf{B}^{>0}$  by writing  $(\mathbf{B}^{>0})^\top \mathbf{B}^{>0} = \widehat{\mathbf{B}}^\top \mathbf{W} \widehat{\mathbf{B}}$ , where  $\widehat{\mathbf{B}}$  is the unweighted 2-commodity edge-vertex incidence matrix with the same edge structure as  $\mathbf{B}^{>0}$ , and  $\mathbf{W}$  is the diagonal matrix of the edge weights. Then, the linear

system in  $\mathcal{MC}_2^{>0}$  can be written as

$$\widehat{\mathbf{B}}^\top \mathbf{W} \widehat{\mathbf{B}} \mathbf{x} = \widehat{\mathbf{B}}^\top \mathbf{W}^{1/2} \mathbf{c}^{\mathbf{B}^{>0}}. \quad (3.41)$$

Before we prove Lemma 3.8.3, we explore some properties of the above linear system.

**Claim 3.8.4.** *There exist a 2-commodity edge-vertex incidence matrix  $\widetilde{\mathbf{B}}$ , a diagonal matrix  $\widetilde{\mathbf{W}}$ , and a vector  $\widetilde{\mathbf{c}}$  such that*

1.  $\widetilde{\mathbf{B}}^\top \widetilde{\mathbf{W}} \widetilde{\mathbf{B}} = \widehat{\mathbf{B}}^\top \mathbf{W} \widehat{\mathbf{B}}$  and  $\widehat{\mathbf{B}}^\top \mathbf{W}^{1/2} \mathbf{c}^{\mathbf{B}^{>0}} = \widetilde{\mathbf{B}}^\top \mathbf{W}^{1/2} \widetilde{\mathbf{c}}$ ,
2. *all the rows of  $\widetilde{\mathbf{B}}$  are distinct.*

*Proof.* If all rows of  $\widehat{\mathbf{B}}$  are distinct, then we set  $\widetilde{\mathbf{B}} = \widehat{\mathbf{B}}$ ,  $\widetilde{\mathbf{W}} = \mathbf{W}$  and  $\widetilde{\mathbf{c}} = \mathbf{c}^{\mathbf{B}^{>0}}$ . Otherwise, assume the  $i$ th row and the  $j$ th row of  $\widehat{\mathbf{B}}$  are the same. Note

$$\widehat{\mathbf{B}}^\top \mathbf{W} \widehat{\mathbf{B}} = \sum_{k \neq i, j} \mathbf{W}_{kk} \widehat{\mathbf{B}}_k^\top \widehat{\mathbf{B}}_k + (\mathbf{W}_{ii} + \mathbf{W}_{jj}) \widehat{\mathbf{B}}_i^\top \widehat{\mathbf{B}}_i,$$

and

$$\widehat{\mathbf{B}}^\top \mathbf{W}^{1/2} \mathbf{c}^{\mathbf{B}^{>0}} = \sum_{k \neq i, j} \mathbf{W}_{kk}^{1/2} \mathbf{c}_k^{\mathbf{B}^{>0}} \widehat{\mathbf{B}}_k^\top + (\mathbf{W}_{ii}^{1/2} \mathbf{c}_i^{\mathbf{B}^{>0}} + \mathbf{W}_{jj}^{1/2} \mathbf{c}_j^{\mathbf{B}^{>0}}) \widehat{\mathbf{B}}_i^\top.$$

We can construct  $\widetilde{\mathbf{B}}$  and  $\widetilde{\mathbf{W}}$  by removing the  $j$ th row of  $\widehat{\mathbf{B}}$ ,  $\mathbf{W}$  and set the corresponding edge weight  $\widetilde{\mathbf{W}}_{ii}$  to be  $\mathbf{W}_{ii} + \mathbf{W}_{jj}$ . We construct  $\widetilde{\mathbf{c}}$  by removing the  $j$ th entry of  $\mathbf{c}^{\mathbf{B}^{>0}}$  and set  $\widetilde{\mathbf{c}}_i$  to be  $(\mathbf{W}_{ii}^{1/2} \mathbf{c}_i^{\mathbf{B}^{>0}} + \mathbf{W}_{jj}^{1/2} \mathbf{c}_j^{\mathbf{B}^{>0}}) / (\mathbf{W}_{ii} + \mathbf{W}_{jj})^{1/2}$ . We repeat the above process to merge the identical rows until we get the desired matrices and vector.  $\square$

**Remark.** Let  $\mathbf{x}^*$  be a minimizer of  $\left\| \mathbf{W}^{1/2} \widehat{\mathbf{B}} \mathbf{x} - \mathbf{c}^{\mathbf{B}^{>0}} \right\|_2$ . We know that  $\mathbf{x}^*$  is a solution of the normal equations  $\widehat{\mathbf{B}}^\top \mathbf{W} \widehat{\mathbf{B}} \mathbf{x} = \widehat{\mathbf{B}}^\top \mathbf{W}^{1/2} \mathbf{c}^{\mathbf{B}^{>0}}$ . By Claim 3.8.4,  $\mathbf{x}^*$  is a solution of  $\widetilde{\mathbf{B}}^\top \widetilde{\mathbf{W}} \widetilde{\mathbf{B}} \mathbf{x} = \widetilde{\mathbf{B}}^\top \widetilde{\mathbf{W}}^{1/2} \widetilde{\mathbf{c}}$ , which in turn gives that  $\mathbf{x}^* \in \arg \min_{\mathbf{x}} \left\| \widetilde{\mathbf{W}}^{1/2} \widetilde{\mathbf{B}} \mathbf{x} - \widetilde{\mathbf{c}} \right\|_2$ . Let  $\mathbf{x}$  be an approximate solution such that

$$\left\| \mathbf{x} - \mathbf{x}^* \right\|_{\widehat{\mathbf{B}}^\top \mathbf{W} \widehat{\mathbf{B}}} \leq \epsilon \left\| \widehat{\mathbf{B}}^\top \mathbf{W}^{1/2} \mathbf{c}^{\mathbf{B}^{>0}} \right\|_{(\widehat{\mathbf{B}}^\top \mathbf{W} \widehat{\mathbf{B}})^\dagger}.$$

By Claim 3.8.4, the above equation is equivalent to

$$\|\mathbf{x} - \mathbf{x}^*\|_{\tilde{\mathbf{B}}^\top \tilde{\mathbf{W}} \tilde{\mathbf{B}}} \leq \epsilon \left\| \tilde{\mathbf{B}}^\top \tilde{\mathbf{W}}^{1/2} \tilde{\mathbf{c}} \right\|_{(\tilde{\mathbf{B}}^\top \tilde{\mathbf{W}} \tilde{\mathbf{B}})^\dagger}.$$

It means that  $\mathbf{x}$  is an approximate solution of  $\min_{\mathbf{x}} \left\| \tilde{\mathbf{W}}^{1/2} \tilde{\mathbf{B}} \mathbf{x} - \tilde{\mathbf{c}} \right\|_2$  with  $\epsilon$ -accuracy, vice versa. Thus, the two minimization problems before and after merging identical rows are equivalent.

Without the loss of generality, we assume that all the rows of  $\hat{\mathbf{B}}$  are distinct. We define *the underlying graph* of  $\hat{\mathbf{B}}$  to be the *simple* graph over  $n$  vertices where vertices  $i$  and  $j$  are connected if and only if  $\hat{\mathbf{B}}$  has at least one type of edges between  $i$  and  $j$ . According to the constructions in  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$  (Algorithm 1) and  $\text{REDUCE } \mathcal{MC}_2 \text{ TO } \mathcal{MC}_2^{>0}$  (Algorithm 4), each edge of the underlying graph has exactly 3 types of edges in  $\hat{\mathbf{B}}$  (that is, type 1, type 2 and type 1 + 2), and the edge weights of  $\mathbf{B}^{>0}$  are either at least 1 or equal to a tiny number  $\delta$  assigned in line 1 of Algorithm 4. We call edge weights which are at least 1 as *large* weights.

**Claim 3.8.5.** *For each edge in the underlying graph of  $\hat{\mathbf{B}}$ , exactly one of its type 1, type 2, and type 1 + 2 edges has large weight.*

*Proof.* Note that the large-weight edges all appear in  $\mathbf{B} \in \mathcal{MC}_2$  constructed in Algorithm 1. An edge in  $\mathbf{B}$  is either a gadget edge in  $\mathcal{B}$  or a non-gadget edge in  $\mathcal{A}$ .

1. By Algorithm 2, all edges in a single  $\mathcal{MC}_2$ -gadget are distinct.
2. We show that each gadget edge is distinct from all other edges. By Algorithm 2, in a 2-commodity gadget, each “old” vertex is connected to a “new” vertex, and all new vertices except  $\mathbf{x}_t$  are independent of the new vertices created in all other gadgets. However, in each gadget including  $\mathbf{x}_t$  (as a new vertex or an odd vertex),  $\mathbf{x}_t$  is connected to new vertices which are created in that gadget and do not appear in any other edge outside that gadget.

3. It is possible that two identical non-gadget edges exist when the original linear system instance has redundant or inconsistent constraints. But note all non-gadget edges are type 1 edges, it implies that two identical edges correspond to two identical rows in  $\widehat{\mathbf{B}}$ . By Claim 3.8.4, we can always merge such identical rows so that all non-gadget edges are distinct.

This completes the proof.  $\square$

*Proof of Lemma 3.8.3.* Note the matrices  $\widehat{\mathbf{B}}$ ,  $\mathbf{W}$  and the vector  $\mathbf{c}^{\mathbf{B} > 0}$  are given by the reductions from the general linear system instance  $\mathbf{A}\mathbf{x} = \mathbf{c}^{\mathbf{A}}$ . The goal of proving Lemma 3.8.3 is to determine the edge-vertex incidence matrix  $\mathbf{N}$ , the cost vector  $\mathbf{c}$ , the demand vectors  $\mathbf{d}^1, \mathbf{d}^2$ , the edge capacity vector  $\mathbf{z}$ , and the flows  $\mathbf{y}^1, \mathbf{y}^2$  such that the intermediate linear systems that arise in an interior point method in Equation (3.40) for solving the Min-cost 2-commodity Flow Problem is exactly the above linear system in Equation (3.41).

We first make the coefficient matrices of the two linear systems in Equation (3.40) and (3.41) to be the same, by choosing the matrix  $\mathbf{N}$  and the vectors  $\mathbf{y}^1, \mathbf{y}^2, \mathbf{d}^1, \mathbf{d}^2, \mathbf{z}$  properly. The coefficient matrix of the equality constraints of the 2-commodity linear programming is

$$\mathbf{M} = \begin{pmatrix} \mathbf{N}^\top & \\ & \mathbf{N}^\top \end{pmatrix}. \quad (3.42)$$

We choose  $\mathbf{N}$  to be the edge-vertex incidence matrix of the underlying graph of  $\widehat{\mathbf{B}}$ . By denoting the residue flow amount along an edge as

$$\mathbf{y}^r \stackrel{\text{def}}{=} \mathbf{z} - \mathbf{y}^1 - \mathbf{y}^2,$$

we can rewrite the Hessian as:

$$\mathbf{H} = \frac{1}{t} \sum_{i \in [m]} \frac{1}{(\mathbf{y}_i^1)^2} \mathbf{e}_{2i-1} \mathbf{e}_{2i-1}^\top + \frac{1}{(\mathbf{y}_i^2)^2} \mathbf{e}_{2i} \mathbf{e}_{2i}^\top + \frac{1}{(\mathbf{y}_i^r)^2} (\mathbf{e}_{2i-1} + \mathbf{e}_{2i}) (\mathbf{e}_{2i-1} + \mathbf{e}_{2i})^\top,$$

where  $\mathbf{e}_i \in \mathbb{R}^{2m}$  is the standard basis vector.

We can rearrange the rows and columns of  $\mathbf{H}$  such that, row  $2i - 1$  and column  $2i - 1$  correspond to the  $\mathbf{y}^1$ -flow of the  $i$ th edge, and row  $2i$  and column  $2i$  correspond to the  $\mathbf{y}^2$ -flow of the  $i$ th edge. After this rearrangement,  $\mathbf{H}$  becomes a block diagonal matrix, where each edge  $i$  corresponds to a one  $2 \times 2$  block given by

$$\mathbf{H}_{2i-1:2i} \stackrel{\text{def}}{=} \begin{pmatrix} \frac{1}{(\mathbf{y}_i^1)^2} + \frac{1}{(\mathbf{y}_i^r)^2} & \frac{1}{(\mathbf{y}_i^r)^2} \\ \frac{1}{(\mathbf{y}_i^r)^2} & \frac{1}{(\mathbf{y}_i^2)^2} + \frac{1}{(\mathbf{y}_i^r)^2} \end{pmatrix}. \quad (3.43)$$

The block diagonal structure of  $\mathbf{H}$  then gives:

$$\mathbf{H}^{-1} = \text{DIAG} \left( \mathbf{H}_{1:2}^{-1}, \dots, \mathbf{H}_{2m-1:2m}^{-1} \right),$$

and for each edge  $i$  we have:

$$\mathbf{H}_{2i-1:2i}^{-1} = \frac{1}{\alpha_i} \left( \begin{pmatrix} (\mathbf{y}_i^1)^2 (\mathbf{y}_i^r)^2 & 0 \\ 0 & (\mathbf{y}_i^2)^2 (\mathbf{y}_i^r)^2 \end{pmatrix} + (\mathbf{y}_i^1)^2 (\mathbf{y}_i^2)^2 \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} \right), \quad (3.44)$$

where  $\alpha_i = (\mathbf{y}_i^1)^2 + (\mathbf{y}_i^2)^2 + (\mathbf{y}_i^r)^2$ . It means that

1. the type 1 edge has weight  $(\mathbf{y}_i^1)^2 (\mathbf{y}_i^r)^2 / \alpha_i$ ,
2. the type 2 edge has weight  $(\mathbf{y}_i^2)^2 (\mathbf{y}_i^r)^2 / \alpha_i$ , and
3. the type 1 + 2 edge has weight  $(\mathbf{y}_i^1)^2 (\mathbf{y}_i^2)^2 / \alpha_i$ .

Note these three types of weights are symmetric. Let  $w_{(i,1)}, w_{(i,2)}, w_{(i,1+2)}$  be the edge weights of type 1, type 2 and type 1 + 2 edge of the  $i$ th edge, respectively. These edge

weights are given by the corresponding diagonals of  $\mathbf{W}$ . We set the values of  $\mathbf{y}_i^1, \mathbf{y}_i^2, \mathbf{y}_i^r$  such that

$$\begin{aligned}(\mathbf{y}_i^1)^2(\mathbf{y}_i^r)^2/\alpha_i &= w_{(i,1)}, \\(\mathbf{y}_i^2)^2(\mathbf{y}_i^r)^2/\alpha_i &= w_{(i,2)}, \\(\mathbf{y}_i^1)^2(\mathbf{y}_i^2)^2/\alpha_i &= w_{(i,1+2)}.\end{aligned}$$

Solving the above equations gives:

$$\begin{aligned}(\mathbf{y}_i^1)^2 &= \frac{w_{(i,1)}w_{(i,1+2)}}{w_{(i,2)}} + w_{(i,1)} + w_{(i,1+2)}, \\(\mathbf{y}_i^2)^2 &= \frac{w_{(i,2)}w_{(i,1+2)}}{w_{(i,1)}} + w_{(i,2)} + w_{(i,1+2)}, \\(\mathbf{y}_i^r)^2 &= \frac{w_{(i,1)}w_{(i,2)}}{w_{(i,1+2)}} + w_{(i,1)} + w_{(i,2)}.\end{aligned}\tag{3.45}$$

The  $\mathbf{y}^1, \mathbf{y}^2$  and  $\mathbf{y}^r$  determine the edge-capacity vector  $\mathbf{z}$ , and they together with  $\mathbf{N}$  determine the vertex-demand vectors  $\mathbf{d}^1$  and  $\mathbf{d}^2$ .

Since we rearranged the rows and columns of  $\mathbf{H}$ , we need do the same rearrangement for  $\mathbf{M}$ . After rearranging row and columns of  $\mathbf{M}$ ,  $\mathbf{M}$  is of the form that, column  $2i - 1$  and column  $2i$  correspond to the  $\mathbf{y}^1$ -flow and  $\mathbf{y}^2$ -flow of the  $i$ th edge, respectively, and row  $2i - 1$  and row  $2i$  correspond to the  $\mathbf{u}$ -coordinate and  $\mathbf{v}$ -coordinate of vertex  $i$ , respectively. In  $\mathbf{M}$ , each edge has a  $2 \times 2$  identity matrix for one endpoint, and a negative  $2 \times 2$  identity matrix for the other endpoint.

By the above setting, we can check that

$$\mathbf{M}\mathbf{H}^{-1}\mathbf{M}^\top = \widehat{\mathbf{B}}^\top \mathbf{W} \widehat{\mathbf{B}}.\tag{3.46}$$

Note that  $\mathbf{M}$  has dimension  $2n \times 2m$ , and  $\widehat{\mathbf{B}}$  has dimension  $3m \times 2n$ .

Secondly, we make the right hand side vectors of the two linear systems in Equa-

tion (3.40) and (3.41) to be the same, that is,

$$-t\mathbf{M}\mathbf{H}^{-1}\mathbf{g} = \widehat{\mathbf{B}}^\top \mathbf{W}^{1/2} \mathbf{c}^{\mathbf{B}^{>0}}, \quad (3.47)$$

by choosing the cost vector  $\mathbf{c}$  properly. The gradient  $\mathbf{g}$  is given by:

$$\mathbf{g} = \mathbf{c} - \frac{1}{t} \sum_{i \in [m]} \frac{1}{\mathbf{y}_i^1} \mathbf{e}_{2i-1} + \frac{1}{\mathbf{y}_i^2} \mathbf{e}_{2i} - \frac{1}{\mathbf{y}_i^r} (\mathbf{e}_{2i-1} + \mathbf{e}_{2i}).$$

For simplicity, let

$$\mathbf{f}(\mathbf{y}) \stackrel{\text{def}}{=} t(\mathbf{c} - \mathbf{g}).$$

$\mathbf{f}(\mathbf{y})$  is a vector in  $2m$ -dimension such that, the  $i$ th edge corresponds to the  $2 \times 1$  sub-vector

$$\begin{pmatrix} \frac{1}{\mathbf{y}_i^1} - \frac{1}{\mathbf{y}_i^r} \\ \frac{1}{\mathbf{y}_i^2} - \frac{1}{\mathbf{y}_i^r} \end{pmatrix}.$$

Then we have

$$-t\mathbf{M}\mathbf{H}^{-1}\mathbf{g} = \mathbf{M}\mathbf{H}^{-1}(\mathbf{f} - t\mathbf{c}).$$

Recall the right hand side vector of Equation (3.47) is  $\widehat{\mathbf{B}}^\top \mathbf{W}^{1/2} \mathbf{c}^{\mathbf{B}^{>0}}$ . According to our construction in  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$  (Algorithm 1) and  $\text{REDUCE } \mathcal{MC}_2 \text{ TO } \mathcal{MC}_2^{>0}$  (Algorithm 4),  $\mathbf{c}^{\mathbf{B}^{>0}} = (\mathbf{c}^{\mathbf{A}}; \mathbf{0})$ , where  $\mathbf{c}^{\mathbf{A}}$  is the right hand side vector of the general linear system instance. Note all the nonzero entries of  $\mathbf{c}^{\mathbf{B}^{>0}}$  correspond to the type 1 edges in the main constraint set  $\mathcal{A}$ , see line 16 of Algorithm 1.

Recall the structure of  $\mathbf{M}$  in Equation (3.42), the edge-vertex incidence matrix  $\mathbf{N}$  contains all the  $m$  edges of the underlying graphs, which are exactly the set of edges in  $\mathbf{B}$  constructed in  $\text{REDUCE } \mathcal{G}_{z,2} \text{ TO } \mathcal{MC}_2$  (Algorithm 1). Note for each edge of  $\mathbf{N}$ , its large-weight edge can be any of the corresponding type 1, type 2 and type 1 + 2 edges. By Claim 3.8.5, we can partition the underlying graph edges into  $E_1 \cup E_2 \cup E_{1+2}$  such that  $E_k$

contains all the underlying edges whose large-weight edges are the corresponding type  $k$  edges,  $k \in \{1, 2, 1 + 2\}$ . Note  $\mathbf{c}_{(i,1)}^{\mathbf{B}^{>0}} \neq 0$  only if edge  $i$  is in  $E_1$ .

We define a vector  $\mathbf{c}' \in \mathbb{R}^{2m}$  such that the entry of  $\mathbf{c}'$  which corresponds to the  $\mathbf{y}^1$ -flow of the  $i$ th edge equals to  $w_{(i,1)}^{1/2} \mathbf{c}_{(i,1)}^{\mathbf{B}^{>0}}$ . This definition gives:

$$\widehat{\mathbf{B}}^\top \mathbf{W}^{1/2} \mathbf{c}^{\mathbf{B}^{>0}} = \mathbf{M} \mathbf{c}'. \quad (3.48)$$

By Equation (3.47),

$$\mathbf{M} \mathbf{H}^{-1} (\mathbf{f} - t\mathbf{c}) = \mathbf{M} \mathbf{c}'.$$

Rearranging it,

$$\mathbf{M} \mathbf{H}^{-1} (\mathbf{f} - t\mathbf{c} - \mathbf{H} \mathbf{c}') = \mathbf{0}.$$

According to  $\mathcal{MC}_2$ GADGET (Algorithm 2), each gadget has a cycle containing both the 2 edges in  $E_{1+2}$ :  $(t + 3, t + 4, t + 2, t + 5, t + 6, t + 1, t + 3)$ . Let  $E'$  be the set of all the edges in such gadget cycles. Let  $\mathbf{p} \in \{0, 1\}^m$  such that each entry corresponding to an edge in  $E'$  has value 1. Let  $\mathbf{q} \stackrel{\text{def}}{=} (\mathbf{p}; \mathbf{p})$ . According to REDUCE  $\mathcal{G}_{z,2}$  TO  $\mathcal{MC}_2$  (Algorithm 1), every gadget edge  $i$  has  $\mathbf{c}'_i = 0$ . Thus, we are free to choose the sign of a column of  $\mathbf{N}^\top$  which corresponds to a gadget edge, which does not change Equations (3.46) and (3.48).

<sup>3</sup> Without loss of generality, we assume that after proper sign flipping,  $\mathbf{M} \mathbf{q} = \mathbf{0}$ . It gives that

$$\mathbf{M} \mathbf{H}^{-1} (\mathbf{f} - t\mathbf{c} - \mathbf{H} \mathbf{c}' + \mathbf{H} \mathbf{q}) = \mathbf{0}.$$

We set  $t\mathbf{c} = \mathbf{f} - \mathbf{H} \mathbf{c}' + \lambda \mathbf{H} \mathbf{q}$ , where  $\lambda > 0$  is a sufficiently large constant to be determined later.

For each edge  $i$ , let  $\mathcal{I}_i$  be an indicator whose value is 1 if edge  $i$  in the set  $E'$  and 0

---

<sup>3</sup>Note the demand vectors  $\mathbf{d}^1$  and  $\mathbf{d}^2$  change after we change  $\mathbf{N}$ .



otherwise. For the  $i$ th edge, its corresponding  $2 \times 1$  sub-vector in  $t\mathbf{c}$  is

$$t \begin{pmatrix} \mathbf{c}_i \\ \mathbf{c}_{m+i} \end{pmatrix} = \begin{pmatrix} \frac{1}{\mathbf{y}_i^1} - \frac{1}{\mathbf{y}_i^r} \\ \frac{1}{\mathbf{y}_i^2} - \frac{1}{\mathbf{y}_i^r} \end{pmatrix} - \mathbf{c}'_i \begin{pmatrix} \frac{1}{(\mathbf{y}_i^1)^2} + \frac{1}{(\mathbf{y}_i^r)^2} \\ \frac{1}{(\mathbf{y}_i^r)^2} \end{pmatrix} + \mathcal{I}_i \lambda \mathbf{s}^i,$$

where

$$\mathbf{s}^i \stackrel{\text{def}}{=} \begin{pmatrix} \frac{1}{(\mathbf{y}_i^1)^2} + \frac{2}{(\mathbf{y}_i^r)^2} \\ \frac{1}{(\mathbf{y}_i^2)^2} + \frac{2}{(\mathbf{y}_i^r)^2} \end{pmatrix}.$$

Note that  $\mathbf{s}^i \geq \mathbf{0}$  (entry-wise). Plugging the solution of  $\mathbf{y}$  from Equation (3.45) gives:

$$t \begin{pmatrix} \mathbf{c}_i \\ \mathbf{c}_{m+i} \end{pmatrix} = \begin{pmatrix} \sqrt{\frac{w_{(i,2)}}{\gamma}} - \sqrt{\frac{w_{(i,1+2)}}{\gamma}} \\ \sqrt{\frac{w_{(i,1)}}{\gamma}} - \sqrt{\frac{w_{(i,1+2)}}{\gamma}} \end{pmatrix} - \mathbf{c}'_i \begin{pmatrix} \frac{w_{(i,2)} + w_{(i,1+2)}}{\gamma} \\ \frac{w_{(i,1+2)}}{\gamma} \end{pmatrix} + \mathcal{I}_i \lambda \mathbf{s}^i,$$

where  $\gamma \stackrel{\text{def}}{=} w_{(i,1)}w_{(i,1+2)} + w_{(i,1)}w_{(i,2)} + w_{(i,2)}w_{(i,1+2)}$ .

By Claim 3.8.5, each edge of the underlying graph has exactly one type of edges with large weight. We deal with each of these cases separately under the condition that two of the edge's weights are much smaller than that of the third one. We will also denote this ratio using  $\epsilon > 0$ :

1. If  $w_{(i,2)} = w_{(i,1+2)} = \epsilon w_{(i,1)}$ , then

$$t \begin{pmatrix} \mathbf{c}_i \\ \mathbf{c}_{m+i} \end{pmatrix} \geq -\mathbf{c}'_i \begin{pmatrix} \frac{2\epsilon}{\mu} \\ \frac{\epsilon}{\mu} \end{pmatrix} + \begin{pmatrix} 0 \\ \frac{1-\sqrt{\epsilon}}{\sqrt{\mu}} \end{pmatrix},$$

where  $\mu = (2\epsilon + \epsilon^2)w_{(i,1)}$ . This corresponds to an edge of  $E_1$ . If  $\mathbf{c}'_i \leq 0$ , then  $t\mathbf{c}_i, t\mathbf{c}_{m+i} \geq 0$ . If  $\mathbf{c}'_i > 0$ , then  $\mathbf{c}'_i = w_{(i,1)}^{1/2} \mathbf{c}_{(i,1)}^{\text{B} > 0}$  where  $\mathbf{c}_{(i,1)}^{\text{B} > 0}$  comes from an entry of  $\mathbf{c}^A$ . We multiply -1 on both sides of the corresponding equation in the general linear system instance  $\mathbf{A}\mathbf{x} = \mathbf{c}^A$  which creates this edge  $i$ , and will get  $\mathbf{c}'_i < 0$  instead.<sup>4</sup>

---

<sup>4</sup>We first flip signs of equations of  $\mathbf{A}\mathbf{x} = \mathbf{c}^A$  so that  $\mathbf{c}^A \leq \mathbf{0}$  (entry-wise), which determines  $\mathbf{N}$ ,  $\mathbf{M}$  and  $\mathbf{c}'$ . We then flip signs of columns of  $\mathbf{N}$ , which we guarantee does not violate Equation (3.46) and (3.48).

2. If  $w_{(i,1)} = w_{(i,1+2)} = \epsilon w_{(i,2)}$ , then

$$t \begin{pmatrix} \mathbf{c}_i \\ \mathbf{c}_{m+i} \end{pmatrix} \geq -\mathbf{c}'_i \begin{pmatrix} \frac{1+\epsilon}{\mu} \\ \frac{\epsilon}{\mu} \end{pmatrix} + \begin{pmatrix} \frac{1-\sqrt{\epsilon}}{\sqrt{\mu}} \\ 0 \end{pmatrix},$$

where  $\mu = (2\epsilon + \epsilon^2)w_{(i,2)}$ . This corresponds to an edge of  $E_2$ , in which  $\mathbf{c}'_i = 0$ . We can check that  $t\mathbf{c}_i, t\mathbf{c}_{m+i} \geq 0$ .

3. If  $w_{(i,1)} = w_{(i,2)} = \epsilon w_{(i,1+2)}$ , then

$$t \begin{pmatrix} \mathbf{c}_i \\ \mathbf{c}_{m+i} \end{pmatrix} = -\mathbf{c}'_i \begin{pmatrix} \frac{1+\epsilon}{\mu} \\ \frac{1}{\mu} \end{pmatrix} + \begin{pmatrix} \frac{\sqrt{\epsilon}-1}{\sqrt{\mu}} \\ \frac{\sqrt{\epsilon}-1}{\sqrt{\mu}} \end{pmatrix} + \mathcal{I}_i \lambda \mathbf{s}^i,$$

where  $\mu = (2\epsilon + \epsilon^2)w_{(i,1+2)}$ . This corresponds to an edge of  $E_{1+2}$ , in which  $\mathbf{c}'_i = 0$  and  $\mathcal{I}_i = 1$ . We choose  $\lambda$  such that, for every edge  $i$  in this case,  $t\mathbf{c}_i, t\mathbf{c}_{m+i} \geq 0$ .

□

Note in IPMs, the values of  $t$  is a sequence of increasing values. Whenever  $t \geq m/\epsilon'$ , the optimality gap is smaller than  $\epsilon'$ . This ensures that the cost vector  $\mathbf{c}$ 's entries are in a reasonable range.

### 3.8.3 Isotropic Total Variation Minimization

For the isotropic total variation problem, we follow the formulations given in [32], namely given a graph  $G = (V, E, \mathbf{w})$ , we partition the sets into  $S_1 \dots S_k$ , and minimize the objective

$$\|\mathbf{y} - \mathbf{s}\|_2^2 + \sum_{1 \leq i \leq k} \mathbf{w}_i \sqrt{\sum_{(u,v) \in S_i} \mathbf{w}_{uv} (\mathbf{y}_u - \mathbf{y}_v)^2}.$$

where  $\mathbf{s}$  is the input signal vector. Let  $\mathbf{N}$  be the edge-vertex incidence matrix of the graph  $G$ , then the dual of the grouped flow problem is:

$$\begin{aligned} \max \quad & \mathbf{s}^\top \mathbf{N}^\top \mathbf{f} \\ \text{subject to:} \quad & \sum_{e \in S_i} \mathbf{w}_e^{-1} \mathbf{f}_e^2 \leq \mathbf{w}_i^{-1} \quad \forall 1 \leq i \leq k \end{aligned}$$

Both the primal and dual problems can be formulated into log barriers via second order cones [76]. In the case with minimizing vertex labels, we introduce a variable  $\mathbf{y}_i$  for each cluster  $i$ , and instead minimize  $\sum_i \mathbf{w}_i \mathbf{y}_i$  subject to the constraint

$$\mathbf{y}_i^2 \geq \sum_{(u,v) \in S_i} \mathbf{w}_e^{-1} (\mathbf{y}_u - \mathbf{y}_v)^2,$$

which in turn leads to the logarithmic barrier function

$$\log \left( \mathbf{y}_i^2 - \sum_{(u,v) \in S_i} \mathbf{w}_e^{-1} (\mathbf{y}_u - \mathbf{y}_v)^2 \right),$$

while in the flow case the barrier functions for each set of edges  $S_i$  is

$$\log \left( \mathbf{w}_i^{-1} - \sum_{e \in S_i} \mathbf{w}_e^{-1} \mathbf{f}_e^2 \right).$$

Due to the connections with the multicommodity flow problems, we will only state the connections through the flow version here. Recall that linear systems in  $\mathcal{MC}_2$  have the form:

$$\mathbf{L}^1 \otimes \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + \mathbf{L}^2 \otimes \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} + \mathbf{L}^{1+2} \otimes \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix},$$

where  $E_1, E_2, E_{1+2}$  denotes the edges in  $\mathbf{L}^1, \mathbf{L}^2, \mathbf{L}^{1+2}$ , respectively. Edges in  $E_1 \cup E_2$

correspond to individual edges, so it remains to show that edges in  $E_{1+2}$  can correspond to the Hessian matrices of clusters containing pairs of edges. In a cluster with two edges  $\mathbf{f}_1$  and  $\mathbf{f}_2$ , the gradient of the function

$$\log(\alpha - \mathbf{w}_1^{-1}\mathbf{f}_1^2 - \mathbf{w}_2^{-1}\mathbf{f}_2^2)$$

is

$$\frac{2}{r} \begin{pmatrix} -\mathbf{w}_1^{-1}\mathbf{f}_1 \\ -\mathbf{w}_2^{-1}\mathbf{f}_2 \end{pmatrix},$$

where we treat

$$r \stackrel{\text{def}}{=} \alpha - \mathbf{w}_1^{-1}\mathbf{f}_1^2 - \mathbf{w}_2^{-1}\mathbf{f}_2^2$$

as a new free variable because we are free to choose  $\alpha$  in the IPM instance. Differentiating again (via the product rule) then gives the Hessian matrix:

$$\frac{-4}{r^2} \begin{pmatrix} \mathbf{w}_1^{-2}\mathbf{f}_1^2 & \mathbf{w}_1^{-1}\mathbf{w}_2^{-2}\mathbf{f}_1\mathbf{f}_2 \\ \mathbf{w}_1^{-1}\mathbf{w}_2^{-2}\mathbf{f}_1\mathbf{f}_2 & \mathbf{w}_2^{-2}\mathbf{f}_2^2 \end{pmatrix} + \frac{-2}{r} \begin{pmatrix} \mathbf{w}_1^{-1} & 0 \\ 0 & \mathbf{w}_2^{-1} \end{pmatrix}.$$

As we are free to choose  $\mathbf{w}_1$  and  $\mathbf{w}_2$ , the second matrix can be any diagonal matrix. Then the freedom in choosing  $\mathbf{f}_1$  and  $\mathbf{f}_2$  then means the first matrix can be any rank 1 object. Therefore, this Hessian is equivalent to the block Hessian for two commodity flows generated in Equation (3.43) (up to a change of sign, since we are dealing with a maximization problem here), and its inverse as given in Equation (3.44) provides the characterization from Definition 2.3.4. We can also check more directly that the  $1 + 2$  edge can be represented as one of these Hessian inverse blocks.

**Claim 3.8.6.** *For each edge  $(i, j)$  in  $E_{1+2}$ , there exist an edge-vertex incidence matrix  $\mathbf{N}$ ,*

a diagonal matrix  $\mathbf{W}$  and a vector  $\mathbf{r}$  such that  $\mathbf{W} \succcurlyeq \mathbf{r}\mathbf{r}^\top$  and

$$\mathbf{L}_{ij}^{1+2} \otimes \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} = \mathbf{N}^\top (\mathbf{W} - \mathbf{r}\mathbf{r}^\top) \mathbf{N}.$$

*Proof.* For simplicity, we remove all zero rows and columns of  $\mathbf{L}_{ij}^{1+2}$  so that the 2-commodity matrix for edge  $(i, j)$  only has dimension  $4 \times 4$ . We pick

$$\mathbf{N} = \begin{pmatrix} 1 & 0 \\ -1 & 0 \\ 0 & 1 \\ 0 & -1 \end{pmatrix}, \mathbf{W} = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \text{ and } \mathbf{r} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

We check the PSD condition,

$$\mathbf{W} - \mathbf{r}\mathbf{r}^\top = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} \succcurlyeq \mathbf{0}.$$

Then, we check that this decomposition equals to the 2-commodity matrix

$$\begin{aligned} \mathbf{N}^\top (\mathbf{W} - \mathbf{r}\mathbf{r}^\top) \mathbf{N} &= \begin{pmatrix} 1 & 0 \\ -1 & 0 \\ 0 & 1 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix} \\ &= \begin{pmatrix} 1 \\ -1 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 & -1 & -1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & -1 & -1 & 1 \\ -1 & 1 & 1 & -1 \\ 1 & -1 & -1 & 1 \\ -1 & 1 & 1 & -1 \end{pmatrix} = \mathbf{L}_{ij}^{1+2} \otimes \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}. \end{aligned}$$

□

*Proof of Lemma 3.2.9.* Since the linear system related to the isotropic total variation minimization matrix is the same as the linear system for 2-commodity, all complexity parameters of these two linear systems are the same.  $\square$

## CHAPTER 4

### HARDNESS RESULTS OF PACKING LPS

In this chapter, we prove Theorem 1.1.2 and Theorem 1.1.3 by proving the following theorems.

**Theorem 4.0.7.** *If one can  $(1 - \epsilon)$ -approximately solve a packing LP of dimension  $n \times n$  in time  $O(n^\beta / \epsilon^\alpha)$  for some constant  $\alpha, \beta > 0$ , then one can solve any linear system instance  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon)$  of dimension  $O(n) \times O(n)$  in time  $\tilde{O}(n^{\beta+2\alpha} \kappa^{2\alpha}(\mathbf{A}) \log(1/\epsilon))$ , where  $\kappa(\mathbf{A})$  is the condition number.*

Considering  $\beta = 2$  and  $\alpha = 0.0009$ , Theorem 4.0.7 states that:  $(1 - \epsilon)$ -approximately solving a packing LP of dimension  $n \times n$  in time  $O(n^2 / \epsilon^{0.0009})$  would imply approximately solving a linear system of dimension  $n \times n$  and condition number  $n^{10}$  in time  $O(n^{2.1} \log(1/\epsilon))$ . Given the best known bound on the matrix multiplication  $\omega < 2.373$ , this would prove Theorem 1.1.2.

**Theorem 4.0.8.** *If one can  $(1 - \epsilon)$ -approximately solve a packing LP with  $N$  non-zeros in time  $O(N / \epsilon^\alpha)$  for some constant  $\alpha > 0$ , then one can solve any linear system instance  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon)$  in time  $\tilde{O}(\text{nnz}(\mathbf{A})^{1+1.5\alpha} \kappa^{2\alpha}(\mathbf{A}) \log(1/\epsilon))$ , where  $\kappa(\mathbf{A})$  is the condition number.*

Consider  $\alpha = 0.22$ , Theorem 4.0.8 states that:  $(1 - \epsilon)$ -approximately solving a packing LP with  $N$  non-zeros in time  $O(N / \epsilon^{0.22})$  would imply approximately solving a linear system with  $N$  non-zeros and condition number  $\kappa = n^{1.5}$  in time  $o(N \kappa \log(1/\epsilon))$ . This is asymptotically faster than Conjugate Gradient, and thus proves Theorem 1.1.3. Note here we do not require  $\mathbf{A}$  to be symmetric and PSD, Conjugate Gradient in fact solves  $\mathbf{A}^\top \mathbf{A}$ . Since  $\sqrt{\kappa(\mathbf{A}^\top \mathbf{A})} = \kappa(\mathbf{A})$ , the above running time is consistent to the usually used running time of CG.

We remark that for simplicity our proofs only consider sequential solvers. We believe that the same ideas work for parallel solvers.

#### 4.1 Reducing a Linear System to a Packing LP

In this section we discuss how to reduce solving linear systems to solving packing LPs. We introduce two reductions from linear equation problems to packing LP problems. Using the two reductions, we will prove Theorems 4.0.7 and 4.0.8 later in this section.

**Normalizing a Linear System.** Given a linear system instance  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon)$ , where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , we denote  $\mathbf{p} \stackrel{\text{def}}{=} \mathbf{A}^\top \mathbf{b}$ . WLOG we assume  $\mathbf{A}, \mathbf{b}$  are normalized so that

$$\|\mathbf{A}\|_{\max} \stackrel{\text{def}}{=} \max_{i,j} |\mathbf{A}_{i,j}| = 1, \quad \|\mathbf{p}\|_2 = 1.$$

This normalization can be done in  $\text{nnz}(\mathbf{A})$  time. Given any  $\epsilon$ -approximate solution  $\mathbf{x}$  to the normalized linear system, that is,  $\|\mathbf{Ax} - \Pi_A \mathbf{b}\|_2 \leq \epsilon \|\Pi_A \mathbf{b}\|_2$  (refer to Section 2.2), we can scale the solution  $\mathbf{x}$  to get an  $\epsilon$ -approximate solution to the original linear system. Moreover, the condition number  $\kappa(\mathbf{A})$  is not changed by the normalization. Furthermore, as  $\|\mathbf{A}\|_{\max} = 1$ , we must have  $\sigma_{\max}(\mathbf{A}) \geq 1$ , which gives the following claim.

**Claim 4.1.1.**  $\sigma_{\min}^{-1}(\mathbf{A}) \leq \kappa(\mathbf{A})$ .

Note that  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon)$  does not guarantee that  $\mathbf{b}$  is in the column space of  $\mathbf{A}$ , which means that  $\mathbf{Ax} = \mathbf{b}$  is infeasible. As suggested in Section 2.2, we rewrite  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon)$  as the following:

$$\begin{aligned} \mathbf{A}^\top \mathbf{z} &= \mathbf{p} \\ \mathbf{Ax} &= \mathbf{z} \end{aligned} \tag{4.1}$$

Here, recall that  $\mathbf{p} = \mathbf{A}^\top \mathbf{b}$ , and  $\mathbf{z}, \mathbf{x}$  are unknown vector variables. Note Equation (4.1)



always has a solution:

$$\mathbf{z} = (\mathbf{A}\mathbf{A}^\top)^\dagger \mathbf{A}\mathbf{A}^\top \mathbf{b}, \quad \mathbf{x} = (\mathbf{A}^\top \mathbf{A})^\dagger \mathbf{A}^\top \mathbf{b}.$$

As explained in Section 2.2,  $\mathbf{x}$  is a minimizer for the regression problem  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$ . In addition, the following claim upper bounds the  $\ell_1$  norm of a solution of Equation (4.1).

**Claim 4.1.2.** *There exists  $\mathbf{x}^*, \mathbf{z}^*$  satisfying the linear system (4.1), and*

$$\|\mathbf{z}^*\|_1 \leq \sqrt{m} \|\Pi_A \mathbf{b}\|_2, \quad \|\mathbf{x}^*\|_1 \leq \sqrt{n} \sigma_{\min}^{-1}(\mathbf{A}) \|\Pi_A \mathbf{b}\|_2,$$

moreover, we know

$$\|\Pi_A \mathbf{b}\|_2 \in \left[ \frac{\sigma_{\min}(\mathbf{A})}{\sigma_{\max}^2(\mathbf{A})} \|\mathbf{p}\|_2, \frac{\sigma_{\max}(\mathbf{A})}{\sigma_{\min}^2(\mathbf{A})} \|\mathbf{p}\|_2 \right].$$

*Proof.* Consider the solution where  $\mathbf{z}^* = \Pi_A \mathbf{b} = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{p}$  and  $\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$ .

It is straightforward to check  $\mathbf{z}^*, \mathbf{x}^*$  satisfy (4.1). Moreover,

$$\|\mathbf{z}^*\|_1 \leq \sqrt{m} \|\mathbf{z}^*\|_2 = \sqrt{m} \|\Pi_A \mathbf{b}\|_2,$$

and

$$\|\mathbf{x}^*\|_1 \leq \sqrt{n} \|\mathbf{x}^*\|_2 \leq \sqrt{n} \sigma_{\min}^{-1}(\mathbf{A}) \|\Pi_A \mathbf{b}\|_2$$

where the last inequality follows from  $\mathbf{A}\mathbf{x}^* = \Pi_A \mathbf{b}$ .

Finally, the upper and lower bounds of  $\|\Pi_A \mathbf{b}\|_2$  follow from  $\Pi_A \mathbf{b} = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{p}$ .

□

**Main Ideas.** Given Equation (4.1), we first apply standard tricks to turn each equality into two inequalities, and make all variables non-negative. For each variable  $x_j$  ( $z_i$ ), we

create two non-negative variables  $\mathbf{x}_j^{(+)}, \mathbf{x}_j^{(-)}$  ( $\mathbf{z}_i^{(+)}, \mathbf{z}_i^{(-)}$ , respectively) such that:

$$\begin{aligned}\mathbf{x}_j &= \mathbf{x}_j^{(+)} - \mathbf{x}_j^{(-)} \\ \mathbf{z}_i &= \mathbf{z}_i^{(+)} - \mathbf{z}_i^{(-)}\end{aligned}$$

Let  $\mathbf{x}^{(+)}, \mathbf{x}^{(-)}, \mathbf{z}^{(+)},$  and  $\mathbf{z}^{(-)}$  be vectors whose entries are  $\mathbf{x}_j^{(+)}, \mathbf{x}_j^{(-)}, \mathbf{z}_i^{(+)}, \mathbf{z}_i^{(-)}$  respectively. We define  $\text{LP}(\mathbf{A}, \mathbf{b})$  as the feasibility LP of finding a solution satisfying:

$$\begin{aligned}\mathbf{A}^T \mathbf{z}^{(+)} - \mathbf{A}^T \mathbf{z}^{(-)} &\leq \mathbf{p} \\ -\mathbf{A}^T \mathbf{z}^{(+)} + \mathbf{A}^T \mathbf{z}^{(-)} &\leq -\mathbf{p} \\ \mathbf{A} \mathbf{x}^{(+)} - \mathbf{A} \mathbf{x}^{(-)} - (\mathbf{z}^{(+)} - \mathbf{z}^{(-)}) &\leq \mathbf{0} \\ -(\mathbf{A} \mathbf{x}^{(+)} - \mathbf{A} \mathbf{x}^{(-)}) + (\mathbf{z}^{(+)} - \mathbf{z}^{(-)}) &\leq \mathbf{0} \\ \mathbf{x}^{(+)}, \mathbf{x}^{(-)}, \mathbf{z}^{(+)}, \mathbf{z}^{(-)} &\geq \mathbf{0}\end{aligned}$$

For simplicity, we use

$$\tilde{\mathbf{A}} \tilde{\mathbf{x}} \leq \tilde{\mathbf{b}}, \tilde{\mathbf{x}} \geq \mathbf{0} \tag{4.2}$$

to denote the above feasibility LP, and explain how to convert all entries of  $\tilde{\mathbf{A}}, \tilde{\mathbf{b}}$  to be non-negative. We will use a *bounding box* constraint:  $\sum_{i=1}^n \tilde{\mathbf{x}}_i \leq U$ , where  $U$  is an upper bound of the  $\ell_1$  norm of some feasible solution  $\tilde{\mathbf{x}}$ . We will show in Claim 4.1.2 that  $U$  is not too large. This is equivalent to the following equality with a new slack variable  $x_0 \geq 0$ :

$$x_0 + \sum_{i=1}^n \tilde{\mathbf{x}}_i = U. \tag{4.3}$$

Adding this equality to each constraint of  $\tilde{\mathbf{A}} \tilde{\mathbf{x}} \leq \tilde{\mathbf{b}}$  makes all entries of the LP non-negative, given that all entries of  $\tilde{\mathbf{A}}, \tilde{\mathbf{b}}$  are between  $[-1, 1]$ . Moreover, this operation does not change the feasible solutions of LP (4.2) under the bounding box constraint.

However, we cannot have the equality (4.3) as a constraint, since it violates the form of packing LPs. Instead, we add the inequality  $\mathbf{x}_0 + \sum_{i=1}^n \tilde{\mathbf{x}}_i \leq U$  as a constraint, and use

$$\max \mathbf{x}_0 + \sum_{i=1}^n \tilde{\mathbf{x}}_i$$

as the objective function. Note there always exists a feasible solution  $\mathbf{x}_0, \tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$  with objective value  $U$ , which forces  $\mathbf{x}_0 + \sum_{i=1}^n \tilde{\mathbf{x}}_i = U$  holds for optimal solutions. Thus, an optimal solution of the above packing LP is a solution for  $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$ .

We will show in Section 4.1.1, the above idea gives a reduction from an arbitrary linear system instance to a packing LP instance such that, the dimensions, and accuracy parameters of the two instances are equivalent up to polynomially factors. In Section 4.1.2, we will present a slightly modified reduction which preserves the sparsity of a linear system instance. Finally, in Section 4.1.3, we will use the standard technique *iterative refinement* for solving linear systems to reduce the polynomially error dependence into logarithmic dependence.

#### 4.1.1 Dense Reduction $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$

$\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$  takes inputs  $\mathbf{A}, \mathbf{b}$  and a tight upper bound of  $\|\mathbf{z}^*\|_1 + \|\mathbf{x}^*\|_1$ , and outputs a packing LP.

To turn  $\text{LP}(\mathbf{A}, \mathbf{b})$  into a packing LP, we add one non-negative slack variable  $s \geq 0$ , and denote

$$\alpha_{\text{sum}} \stackrel{\text{def}}{=} s + \sum_{1 \leq j \leq n} \left( \mathbf{x}_j^{(+)} + \mathbf{x}_j^{(-)} \right) + \sum_{1 \leq i \leq m} \left( \mathbf{z}_i^{(+)} + \mathbf{z}_i^{(-)} \right).$$

Note  $\alpha_{\text{sum}}$  is merely a short-hand for the sum of all the variables instead of a new variable. For each existing constraint of  $\text{LP}(\mathbf{A}, \mathbf{b})$  except the positivity constraints on single variables, we add  $\alpha_{\text{sum}}$  and  $U$  to the LHS and RHS of the constraint respectively, and we also add

$$\alpha_{\text{sum}} \leq U \tag{4.4}$$

as an additional constraint. This completes the construction of  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$ 's constraints, and if written explicitly are the following

$$\begin{aligned}
s + \sum_{1 \leq i \leq m} (1 + \mathbf{A}_{ij}) \mathbf{z}_i^{(+)} + (1 - \mathbf{A}_{ij}) \mathbf{z}_i^{(-)} + \sum_{1 \leq j' \leq n} \mathbf{x}_{j'}^{(+)} + \mathbf{x}_{j'}^{(-)} &\leq \mathbf{p}_j + U & \forall 1 \leq j \leq n \\
s + \sum_{1 \leq i \leq m} (1 - \mathbf{A}_{ij}) \mathbf{z}_i^{(+)} + (1 + \mathbf{A}_{ij}) \mathbf{z}_i^{(-)} + \sum_{1 \leq j' \leq n} \mathbf{x}_{j'}^{(+)} + \mathbf{x}_{j'}^{(-)} &\leq -\mathbf{p}_j + U & \forall 1 \leq j \leq n \\
s + \sum_{1 \leq j \leq n} (1 + \mathbf{A}_{ij}) \mathbf{x}_j^{(+)} + (1 - \mathbf{A}_{ij}) \mathbf{x}_j^{(-)} + 2\mathbf{z}_i^{(-)} + \sum_{1 \leq i' \leq n, i' \neq i} \mathbf{z}_{i'}^{(+)} + \mathbf{z}_{i'}^{(-)} &\leq U & \forall 1 \leq i \leq m \\
s + \sum_{1 \leq j \leq n} (1 - \mathbf{A}_{ij}) \mathbf{x}_j^{(+)} + (1 + \mathbf{A}_{ij}) \mathbf{x}_j^{(-)} + 2\mathbf{z}_i^{(+)} + \sum_{1 \leq i' \leq n, i' \neq i} \mathbf{z}_{i'}^{(+)} + \mathbf{z}_{i'}^{(-)} &\leq U & \forall 1 \leq i \leq m \\
s + \sum_{1 \leq j \leq n} \left( \mathbf{x}_j^{(+)} + \mathbf{x}_j^{(-)} \right) + \sum_{1 \leq i \leq m} \left( \mathbf{z}_i^{(+)} + \mathbf{z}_i^{(-)} \right) &\leq U \\
s, \mathbf{x}^{(+)}, \mathbf{x}^{(-)}, \mathbf{z}^{(+)}, \mathbf{z}^{(-)} &\geq 0
\end{aligned}$$

The objective of  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$  is to maximize  $\alpha_{\text{sum}}$ , that is

$$\max s + \sum_{1 \leq j \leq n} \left( \mathbf{x}_j^{(+)} + \mathbf{x}_j^{(-)} \right) + \sum_{1 \leq i \leq m} \left( \mathbf{z}_i^{(+)} + \mathbf{z}_i^{(-)} \right).$$

**Lemma 4.1.3.** *Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$  constructs a packing LP of size  $O(m + n) \times O(m + n)$ . Furthermore, if  $U \geq \|\mathbf{x}^*\|_1 + \|\mathbf{z}^*\|$  for some feasible solution  $(\mathbf{x}^*, \mathbf{z}^*)$  of  $\text{LP}(\mathbf{A}, \mathbf{b})$ , then  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$  has optimal value  $U$ ; Otherwise,  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$  is infeasible.*

*Proof.* The size of  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$  is obvious. For it to be a packing LP, we need to show

1. All variables are non-negative: This holds by construction.
2. All coefficients are non-negative: This is true since  $\|\mathbf{A}\|_{\max} \leq 1$ , so by adding 1 to every coefficient, all coefficients become non-negative.
3. All constants on the RHS of the constraints are non-negative: This is true since

$\|\mathbf{p}\|_2 = 1$ , so  $\|\mathbf{p}\|_{\infty} \leq 1$ , and we add  $U \geq 1$  to the RHS of each constraint.

4. All coefficients in the objective function are non-negative: This is true since all coefficients are 1 in the objective.

If  $U > \|\mathbf{x}^*\|_1 + \|\mathbf{z}^*\|_1$  for any  $\mathbf{x}^*$  and  $\mathbf{z}^*$  satisfying the linear system (4.1), consider the following solution to  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$  where we set  $\mathbf{x}^{(+)}$  and  $\mathbf{x}^{(-)}$  based on the signs of  $\mathbf{x}^*$

$$\left(\mathbf{x}_i^{(+)}, \mathbf{x}_i^{(-)}\right) := \begin{cases} (\mathbf{x}_i^*, 0) & \text{if } \mathbf{x}_i^* \geq 0, \\ (0, -\mathbf{x}_i^*) & \text{if } \mathbf{x}_i^* \leq 0, \end{cases}$$

and similarly for  $\mathbf{z}^{(+)}$  and  $\mathbf{z}^{(-)}$ :

$$\left(\mathbf{z}_i^{(+)}, \mathbf{z}_i^{(-)}\right) := \begin{cases} (\mathbf{z}_i^*, 0) & \text{if } \mathbf{z}_i^* \geq 0, \\ (0, -\mathbf{z}_i^*) & \text{if } \mathbf{z}_i^* \leq 0, \end{cases}$$

For the added slack variable  $s$ , as  $U > \|\mathbf{x}^*\|_1 + \|\mathbf{z}^*\|_1$ , we can set

$$s := U - \|\mathbf{x}^*\|_1 + \|\mathbf{z}^*\|_1 = U - \left( \sum_{1 \leq j \leq n} \left( \mathbf{x}_j^{(+)} + \mathbf{x}_j^{(-)} \right) + \sum_{1 \leq i \leq m} \left( \mathbf{z}_i^{(+)} + \mathbf{z}_i^{(-)} \right) \right).$$

It is straightforward to see this solution is feasible, and gives objective value  $U$ , which must be the optimal due to the constraint (4.4).  $\square$

In the following lemma, we translate the error bound between the original linear system and  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$ .

**Lemma 4.1.4.** *Consider an instance  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon')$  and any number  $U$  such that  $\text{LP}(\mathbf{A}, \mathbf{x})$  has a feasible solution with  $\ell_1$  norm at most  $U$ . Let  $\epsilon > 0$  be any value satisfying*

$$\epsilon \leq \frac{\epsilon'}{(\sqrt{m} + \sqrt{n})U}.$$

*Suppose we can compute a feasible solution  $\mathbf{x}^{(+)}, \mathbf{x}^{(-)}, \mathbf{z}^{(+)}, \mathbf{z}^{(-)}, \mathbf{s}^{(p)}, \mathbf{s}^{(0)}$  of  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$  with objective value at least  $(1 - \epsilon)U$ . Then,  $\mathbf{x} = \mathbf{x}^{(+)} - \mathbf{x}^{(-)}$  is a solution to  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon')$ ,*

that is,

$$\|\mathbf{A}\mathbf{x} - \Pi_A \mathbf{b}\|_2 \leq \epsilon' \|\Pi_A \mathbf{b}\|_2.$$

*Proof.* To solve  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon')$ , we first check whether  $\mathbf{A}^\top \mathbf{b} = \mathbf{0}$ , that is, whether  $\mathbf{b}$  is orthogonal to the column space of  $\mathbf{A}$ . If there is true, then  $\mathbf{x} = \mathbf{0}$  is a solution to  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon')$  and we have done. In the rest of the proof, we assume that  $\mathbf{A}^\top \mathbf{b} \neq \mathbf{0}$ , equivalently,  $\Pi_A \mathbf{b} \neq \mathbf{0}$ . By Lemma 2.2.5,

$$\|\Pi_A \mathbf{b}\|_2 \geq \frac{1}{\sigma_{\max}(\mathbf{A})} \geq 1. \quad (4.5)$$

Consider any feasible solution, and let  $\mathbf{x} = \mathbf{x}^{(+)} - \mathbf{x}^{(-)}$ ,  $\mathbf{z} = \mathbf{z}^{(+)} - \mathbf{z}^{(-)}$ . We want to bound the error in  $\mathbf{A}\mathbf{x} = \mathbf{z}$ ,  $\mathbf{A}^T \mathbf{z} = \mathbf{p}$ . Consider the  $i$ -th equality constraint in the linear system:  $\mathbf{A}_i \mathbf{x} = \mathbf{z}_i$ , we have a corresponding pair of constraints in  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$

$$\begin{aligned} \alpha_{\text{sum}} - \left( \mathbf{z}_i^{(+)} - \mathbf{z}_i^{(-)} \right) + \sum_{1 \leq j \leq n} \mathbf{A}_{ij} \left( \mathbf{x}_j^{(+)} - \mathbf{x}_j^{(-)} \right) &\leq U \\ \alpha_{\text{sum}} + \left( \mathbf{z}_i^{(+)} - \mathbf{z}_i^{(-)} \right) - \sum_{1 \leq j \leq n} \mathbf{A}_{ij} \left( \mathbf{x}_j^{(+)} - \mathbf{x}_j^{(-)} \right) &\leq U \end{aligned}$$

Since we have a feasible solution, we know from the above two constraints that

$$|\mathbf{A}_i \mathbf{x} - \mathbf{z}_i| \leq U - \alpha_{\text{sum}}.$$

The same argument holds for any linear constraint  $(\mathbf{A}^T)_j \mathbf{z} = \mathbf{p}_j$ , so we have

$$\|\mathbf{A}\mathbf{x} - \mathbf{z}\|_\infty, \|\mathbf{A}^T \mathbf{z} - \mathbf{p}\|_\infty \leq U - \alpha_{\text{sum}},$$

and by the approximation guarantee the RHS is at most  $\epsilon \cdot U$ .

Now we bound the error of the solution  $\mathbf{x}$  in terms of the linear system

$$\begin{aligned}
\|\mathbf{Ax} - \Pi_A \mathbf{b}\|_2 &\leq \|\mathbf{Ax} - \Pi_A \mathbf{z}\|_2 + \|\Pi_A \mathbf{z} - \Pi_A \mathbf{b}\|_2 \\
&\leq \|\mathbf{Ax} - \mathbf{z}\|_2 + \sigma_{\min}^{-1}(A) \|\mathbf{A}^T \mathbf{z} - \mathbf{A}^T \mathbf{b}\|_2 \\
&\leq \sqrt{m} \|\mathbf{Ax} - \mathbf{z}\|_\infty + \sqrt{n} \sigma_{\min}^{-1}(A) \|\mathbf{A}^T \mathbf{z} - \mathbf{A}^T \mathbf{b}\|_\infty \\
&\leq \epsilon \sigma_{\min}^{-1}(A) (\sqrt{m} + \sqrt{n}) U.
\end{aligned} \tag{4.6}$$

In the second inequality, we have  $\|\mathbf{Ax} - \Pi_A \mathbf{z}\|_2 \leq \|\mathbf{Ax} - \mathbf{z}\|_2$  because  $\mathbf{Ax}$  is in the span of columns of  $\mathbf{A}$ , so projecting  $\mathbf{z}$  to the column span of  $\mathbf{A}$  makes the distance to  $\mathbf{Ax}$  smaller. Also in the second inequality, we use

$$\begin{aligned}
\|\Pi_A \mathbf{z} - \Pi_A \mathbf{b}\|_2 &= \|\mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{z} - \mathbf{b})\| \\
&= \sqrt{(\mathbf{z} - \mathbf{b})^T \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{z} - \mathbf{b})} \\
&= \sqrt{(\mathbf{z} - \mathbf{b})^T \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{z} - \mathbf{b})} \\
&\leq \sigma_{\min}^{-1}(A) \|\mathbf{A}^T \mathbf{z} - \mathbf{A}^T \mathbf{b}\|_2.
\end{aligned}$$

Plugging Equation (4.5) into Equation (4.6):

$$\|\mathbf{Ax} - \Pi_A \mathbf{b}\|_2 \leq \epsilon \sigma_{\min}(\mathbf{A})^{-1} (\sqrt{m} + \sqrt{n}) U \cdot \|\Pi_A \mathbf{b}\|_2.$$

Setting

$$\epsilon \leq \frac{\epsilon' \sigma_{\min}(\mathbf{A})}{(\sqrt{m} + \sqrt{n}) U} \leq \frac{\epsilon'}{(\sqrt{m} + \sqrt{n}) U},$$

we have  $\|\mathbf{Ax} - \Pi_A \mathbf{b}\|_2 \leq \epsilon' \|\Pi_A \mathbf{b}\|_2$ . □

By Claim 4.1.2 and Claim 4.1.1, with a  $\log(\kappa(\mathbf{A}))$  overhead in running time, we can assume we have a almost tight bound  $U$ , for example  $U \leq 2\sqrt{m+n}\sigma^{-2}(\mathbf{A})$ . With this  $U$ ,

by Lemma 4.1.4 and Claim 4.1.1, it suffices to set

$$\epsilon = \frac{\epsilon'}{100(m+n)\kappa^2(\mathbf{A})}.$$

#### 4.1.2 Sparse Reduction $\text{PLP}_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$

The sparsity preserving reduction follows the same approach as the dense reduction. However, instead of adding the same  $\alpha_{\text{sum}}$  and  $U$  to the LHS and RHS of every inequality of  $\text{LP}(\mathbf{A}, \mathbf{b})$  to make the coefficients non-negative, we will pin-point only the non-zero coefficients to preserve sparsity. We construct  $\text{PLP}_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$  starting from  $\text{LP}(\mathbf{A}, \mathbf{b})$  as follows.

For each  $j \in [1, n]$ , look at the pair of constraints in  $\text{LP}(\mathbf{A}, \mathbf{b})$  corresponding to the  $j$ -th row of  $\mathbf{A}^T \mathbf{z} = \mathbf{p}$  in the linear system (4.1):

$$\begin{aligned} (\mathbf{A}^T)_j \mathbf{z}^{(+)} - (\mathbf{A}^T)_j \mathbf{z}^{(-)} &\leq \mathbf{p}_j \\ -(\mathbf{A}^T)_j \mathbf{z}^{(+)} + (\mathbf{A}^T)_j \mathbf{z}^{(-)} &\leq -\mathbf{p}_j \end{aligned}$$

We add a non-negative slack variable  $s_j^{(p)}$  that serves as the slack for *both* of these constraints. Furthermore, we denote

$$\alpha_j^{(p)} := s_j^{(p)} + \sum_{i: A_{ij} \neq 0} z_i^{(+)} + z_i^{(-)},$$

and again note  $\alpha_j^{(p)}$  is merely a shorthand rather than a new variable. We add  $\alpha_j^{(p)}$  and  $U$  to the LHS and RHS respectively of both the constraints above. Since  $\|\mathbf{A}\|_{\max} \leq 1$ , all the coefficients in these constraints become non-negative. We then add to  $\text{PLP}_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$  a new inequality

$$\alpha_j^{(p)} \leq U.$$

Similarly, for each  $i \in [1, m]$ , we consider the pair of constraints corresponding to the  $i$ -th



row in  $\mathbf{A}\mathbf{x} = \mathbf{z}$

$$\begin{aligned} \mathbf{A}_i \mathbf{x}^{(+)} - \mathbf{A}_i \mathbf{x}^{(-)} - (\mathbf{z}_i^{(+)} - \mathbf{z}_i^{(-)}) &\leq 0 \\ -(\mathbf{A}_i \mathbf{x}^{(+)} - \mathbf{A}_i \mathbf{x}^{(-)}) + (\mathbf{z}_i^{(+)} - \mathbf{z}_i^{(-)}) &\leq 0 \end{aligned}$$

We add a non-negative slack variable  $s_i^{(0)}$  for both of these constraints, and denote

$$\alpha_i^{(0)} := s_i^{(0)} + \mathbf{z}_i^{(+)} + \mathbf{z}_i^{(-)} + \sum_{1 \leq j \leq n: \mathbf{A}_{ij} \neq 0} \mathbf{x}_j^{(+)} + \mathbf{x}_j^{(-)}.$$

We add  $\alpha_i^{(0)}$  and  $U$  to the LHS and RHS of the pair of constraints, and add  $\alpha_i^{(0)} \leq U$  as a new constraint.

We can write the constraints of  $\text{PLP}_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$  explicitly as follows.

$$\begin{aligned} \mathbf{s}_j^{(p)} + \sum_{1 \leq i \leq m: \mathbf{A}_{ij} \neq 0} (1 + \mathbf{A}_{ij}) \mathbf{z}_i^{(+)} + (1 - \mathbf{A}_{ij}) \mathbf{z}_i^{(-)} &\leq \mathbf{p}_j + U \quad \forall 1 \leq j \leq n \\ \mathbf{s}_j^{(p)} + \sum_{1 \leq i \leq m: \mathbf{A}_{ij} \neq 0} (1 - \mathbf{A}_{ij}) \mathbf{z}_i^{(+)} + (1 + \mathbf{A}_{ij}) \mathbf{z}_i^{(-)} &\leq -\mathbf{p}_j + U \quad \forall 1 \leq j \leq n \\ \mathbf{s}_j^{(p)} + \sum_{1 \leq i \leq m: \mathbf{A}_{ij} \neq 0} \mathbf{z}_i^{(+)} + \mathbf{z}_i^{(-)} &\leq U \quad \forall 1 \leq j \leq n \\ 2\mathbf{z}_i^{(-)} + \mathbf{s}_i^{(0)} + \sum_{1 \leq j \leq n: \mathbf{A}_{ij} \neq 0} (1 + \mathbf{A}_{ij}) \mathbf{x}_j^{(+)} + (1 - \mathbf{A}_{ij}) \mathbf{x}_j^{(-)} &\leq U \quad \forall 1 \leq i \leq m \\ 2\mathbf{z}_i^{(+)} + \mathbf{s}_i^{(0)} + \sum_{1 \leq j \leq n: \mathbf{A}_{ij} \neq 0} (1 - \mathbf{A}_{ij}) \mathbf{x}_j^{(+)} + (1 + \mathbf{A}_{ij}) \mathbf{x}_j^{(-)} &\leq U \quad \forall 1 \leq i \leq m \\ \mathbf{z}_i^{(+)} + \mathbf{z}_i^{(-)} + \mathbf{s}_i^{(0)} + \sum_{1 \leq j \leq n: \mathbf{A}_{ij} \neq 0} \mathbf{x}_j^{(+)} + \mathbf{x}_j^{(-)} &\leq U \quad \forall 1 \leq i \leq m \\ \mathbf{x}^{(+)}, \mathbf{x}^{(-)}, \mathbf{z}^{(+)}, \mathbf{z}^{(-)}, \mathbf{s}^{(p)}, \mathbf{s}^{(0)} &\geq 0 \end{aligned}$$

Finally, the objective of  $\text{PLP}_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$  is to maximize  $\sum_{1 \leq j \leq n} \alpha_j^{(p)} + \sum_{1 \leq i \leq m} \alpha_i^{(0)}$ .

**Lemma 4.1.5.** *Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $\text{PLP}_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$  constructs a packing LP with  $O(\text{nnz}(\mathbf{A}))$  non-zeros. Furthermore, if  $U \geq \|\mathbf{x}^*\|_1 + \|\mathbf{z}^*\|$  for some feasible solution*

$(\mathbf{x}^*, \mathbf{z}^*)$  of  $LP(\mathbf{A}, \mathbf{b})$ , then  $PLP_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$  has optimal value  $(m + n)U$ ; Otherwise,  $PLP_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$  is infeasible.

The proof is a straightforward but tedious adaptation of the proof of Lemma 4.1.3, so we omit it. Now we translate the error bound between the original linear system and  $PLP_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$ .

**Lemma 4.1.6.** *Consider an instance  $LSA(\mathbf{A}, \mathbf{b}, \epsilon')$  and any number  $U$  such that  $LP(\mathbf{A}, \mathbf{x})$  has a feasible solution with  $\ell_1$  norm at most  $U$ . Let  $\epsilon > 0$  be any value satisfying*

$$\epsilon \leq \frac{\epsilon'}{(m + n)U}.$$

*Suppose we can compute a feasible solution  $\mathbf{x}^{(+)}, \mathbf{x}^{(-)}, \mathbf{z}^{(+)}, \mathbf{z}^{(-)}, \mathbf{s}^{(p)}, \mathbf{s}^{(0)}$  of  $PLP_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$  with objective value at least  $(1 - \epsilon) \cdot (m + n)U$ . Then,  $\mathbf{x} = \mathbf{x}^{(+)} - \mathbf{x}^{(-)}$  is a solution to  $LSA(\mathbf{A}, \mathbf{b}, \epsilon')$ , that is,*

$$\|\mathbf{A}\mathbf{x} - \Pi_A \mathbf{b}\|_2 \leq \epsilon' \|\Pi_A \mathbf{b}\|_2.$$

*Proof.* By the same argument as the proof of Lemma 4.1.4, we assume that  $\Pi_A \mathbf{b} \neq \mathbf{0}$  and thus  $\|\Pi_A \mathbf{b}\|_2 \geq \frac{1}{\sigma_{\max}(\mathbf{A})} \geq 1$ .

Given the  $(1 - \epsilon)$ -approx optimal solution to  $PLP_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$ , we consider the following solution to the linear system

$$\mathbf{x} := \mathbf{x}^{(+)} - \mathbf{x}^{(-)}$$

$$\mathbf{z} := \mathbf{z}^{(+)} - \mathbf{z}^{(-)}$$

We want to bound the error in  $\mathbf{A}\mathbf{x} = \mathbf{z}$ ,  $\mathbf{A}^T \mathbf{z} = \mathbf{p}$ .

Consider the  $i$ -th equality constraint in the linear system  $\mathbf{A}_i \mathbf{x} = \mathbf{z}_i$ , we have a corre-

spending pair of constraints in  $\text{PLP}_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$ :

$$\begin{aligned}\alpha_i^{(0)} + \mathbf{A}_i \mathbf{x}^{(+)} - \mathbf{A}_i \mathbf{x}^{(-)} - (\mathbf{z}_i^{(+)} - \mathbf{z}_i^{(-)}) &\leq U \\ \alpha_i^{(0)} - (\mathbf{A}_i \mathbf{x}^{(+)} - \mathbf{A}_i \mathbf{x}^{(-)}) + (\mathbf{z}_i^{(+)} - \mathbf{z}_i^{(-)}) &\leq U\end{aligned}$$

which subtracting away  $\alpha_i^{(0)}$  and  $U$  from the LHS and RHS respectively gives

$$|(\mathbf{A}\mathbf{x})_i - \mathbf{z}_i| \leq U - \alpha_i^{(0)}.$$

The same argument holds for any pair of linear constraints corresponding to  $(\mathbf{A}^T \mathbf{z})_j = \mathbf{p}_j$ , and in total we have

$$\begin{aligned}\|\mathbf{A}\mathbf{x} - \mathbf{z}\|_1 + \|\mathbf{A}^T \mathbf{z} - \mathbf{p}\|_1 \\ \leq (m+n)U - \sum_j \alpha_j^{(p)} - \sum_i \alpha_i^{(0)} \\ \leq \epsilon(m+n)U.\end{aligned}$$

Now we bound the error of the solution  $\mathbf{x}$  in terms of LSA

$$\begin{aligned}\|\mathbf{A}\mathbf{x} - \Pi_A \mathbf{b}\|_2 &\leq \|\mathbf{A}\mathbf{x} - \Pi_A \mathbf{z}\|_2 + \|\Pi_A \mathbf{z} - \Pi_A \mathbf{b}\|_2 \\ &\leq \|\mathbf{A}\mathbf{x} - \mathbf{z}\|_2 + \sigma_{\min}^{-1}(A) \|\mathbf{A}^T \mathbf{z} - \mathbf{A}^T \mathbf{b}\|_2 \\ &\leq \|\mathbf{A}\mathbf{x} - \mathbf{z}\|_1 + \sigma_{\min}^{-1}(A) \|\mathbf{A}^T \mathbf{z} - \mathbf{A}^T \mathbf{b}\|_1 \\ &\leq \epsilon \sigma_{\min}^{-1}(A) (m+n)U \\ &\leq \epsilon \sigma_{\min}^{-1}(A) (m+n)U \cdot \|\Pi_A \mathbf{b}\|_2.\end{aligned}$$

Setting

$$\epsilon \leq \frac{\epsilon'}{(m+n)U},$$

we have  $\|\mathbf{A}\mathbf{x} - \Pi_A \mathbf{b}\|_2 \leq \epsilon' \|\Pi_A \mathbf{b}\|_2$ . □

Similar to the dense reduction, by a standard binary search on  $U$ , it suffices to set

$$\epsilon = \frac{\epsilon'}{100(m^{3/2} + n^{3/2})\kappa^2(\mathbf{A})}.$$

#### 4.1.3 Reducing $\text{poly}(1/\epsilon)$ -Dependence to $\log(1/\epsilon)$ -Dependence

As in Lemma 4.1.4 and Lemma 4.1.6, to solve a linear system instance  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon')$ , it suffices to solve a packing LP constructed by  $\text{PLP}_{\text{dense}}(\mathbf{A}, \mathbf{b}, U)$  or  $\text{PLP}_{\text{sparse}}(\mathbf{A}, \mathbf{b}, U)$  up to accuracy  $\epsilon = \text{poly}(\epsilon, \dim(\mathbf{A})^{-1}, \kappa(\mathbf{A})^{-1})$ . This gives a linear system solver whose running time has  $\text{poly}(1/\epsilon)$ -dependence. We use the standard technique *iterative refinement*, to reduce this  $\text{poly}(1/\epsilon)$ -dependence to  $\log(1/\epsilon)$ -dependence. For completeness, we include a proof below.

**Lemma 4.1.7.** *If we have a solver for  $\text{LSA}(\mathbf{A}, \mathbf{b}, 0.1)$  that works for arbitrary  $\mathbf{b}$ , then we can obtain a solver for  $\text{LSA}(\mathbf{A}, \mathbf{b}, \epsilon)$  by iterating it  $O(\log(1/\epsilon))$  times.*

*Proof.* Consider the linear equation  $\mathbf{A}^\top \mathbf{A} \mathbf{x} = \mathbf{A}^\top \mathbf{b}$ . Let  $\mathbf{M} = \mathbf{A}^\top \mathbf{A}$  and  $\mathbf{p} = \mathbf{A}^\top \mathbf{b}$ . We basically do iterative refinement involving the matrix  $\mathbf{M}$ . Note the desired solution is  $\mathbf{x}^* = \mathbf{M}^{-1} \mathbf{p}$ . (See Fact 2.2.4).

For simplicity let's consider a 2-step version. We start with

$$\mathbf{p}^{(0)} = \mathbf{p},$$

and the LSA solver produces  $\mathbf{x}^{(0)}$  such that

$$\|\mathbf{x}^{(1)} - \mathbf{M}^{-1} \mathbf{p}^{(0)}\|_{\mathbf{M}} \leq 0.1 \|\mathbf{M}^{-1} \mathbf{p}^{(0)}\|_{\mathbf{M}}.$$

Then we define

$$\mathbf{p}^{(1)} := \mathbf{p}^{(0)} - \mathbf{M} \mathbf{x}^{(1)}.$$

Note that at each step, we also have

$$\mathbf{M}\mathbf{x} = \mathbf{A}^T(\mathbf{A}\mathbf{x}),$$

so forming  $\mathbf{b}$  by

$$\mathbf{b}^{(t+1)} = \left( \mathbf{b}^{(t)} - \mathbf{A}\mathbf{x}^{(t)} \right)$$

also works.

By definition we have

$$\mathbf{M}^{-1}\mathbf{p}^{(1)} = \mathbf{M}\mathbf{p}^{(0)} - \mathbf{x}^{(1)},$$

and thus

$$\|\mathbf{M}^{-1}\mathbf{p}^{(1)}\|_{\mathbf{M}} \leq 0.1 \|\mathbf{M}^{-1}\mathbf{p}^{(0)}\|_{\mathbf{M}}.$$

Feeding in  $\mathbf{p}^{(1)}$  as input to the solver in turn gives  $\mathbf{x}^{(2)}$  such that

$$\|\mathbf{x}^{(2)} - \mathbf{M}^{-1}\mathbf{p}^{(1)}\|_{\mathbf{M}} \leq 0.1 \|\mathbf{M}^{-1}\mathbf{p}^{(1)}\|_{\mathbf{M}} \leq 0.01 \|\mathbf{M}^{-1}\mathbf{p}^{(0)}\|_{\mathbf{M}}.$$

Expanding the LHS gives

$$\mathbf{x}^{(2)} - \mathbf{M}^{-1}\mathbf{p}^{(1)} = \mathbf{x}^{(2)} + \mathbf{x}^{(1)} - \mathbf{M}^{-1}\mathbf{p},$$

which means  $\mathbf{x}^{(1)} + \mathbf{x}^{(2)}$  is now a solution with error 0.01. Repeating this gives a  $10\times$  smaller error after each step. □

Combining the Lemmas above, we proved Theorem 4.0.7 and Theorem 4.0.8.

## CHAPTER 5

### 3-D TRUSS LINEAR SYSTEM SOLVER

In this chapter, we prove Theorem 1.1.4 and Theorem 1.1.5.

#### 5.1 Main Algorithm and Proofs of Main Results

The key idea of our algorithm is to combine Nested Dissection and preconditioned iterative algorithms. Motivated by Nested Dissection and the concept  $r$ -division [77], we partition a convex edge-simple 3-D truss (see Definition 2.4.5) into small and separate pieces. Each piece has roughly the same number of vertices, and each piece has a well-shaped boundary consisting of a small number of vertices.

The bottleneck of Nested Dissection is the process of applying Gaussian elimination to the few top-level separators, after eliminating all interior vertices at lower levels (see Section 2.5 for a detailed explanation of nested dissection). At this step, one runs Gaussian elimination on the the Schur complement onto the top-level separators, which is usually dense.

To speed up Nested Dissection, we adapt the idea from support theory. We observe that, after eliminating all interior vertices, the Schur complement onto the boundaries can be well preconditioned by the boundaries alone. Instead of directly running Gaussian elimination on the Schur complement, we run an iterative algorithm Preconditioned Conjugate Gradient (PCG) to solve a sequence of sparse linear systems in the boundaries.

The phenomenon that the boundary itself is a good preconditioner for the Schur complement onto the boundary has a natural interpretation based on structural mechanics. The quadratic form associated with the Schur complement corresponds to the energy associated with deforming the whole truss by squishing or stretching the boundary vertices while leaving the interior intact and finding the positions of the interior vertices that minimize the

overall energy. We show that this energy is not much more than the energy that arises from applying the same deformation to just the boundary tetrahedrons after deleting the interior vertices.

Our main algorithm is presented in Algorithm 12. Its input consists of an edge-simple 3D truss, a column vector, an error parameter, an aspect ratio threshold, and a hollowing rate. Algorithm 12 will use the following subroutines: Algorithm 13, presented in Section 5.4, to compute hollowings of small aspect ratio truss components in line 6; Algorithm 15, presented in Section 5.5, for nested dissection of a partial state of of the preconditioner in line 8.

---

**Algorithm 12** TRUSSSOLVER( $\mathcal{T} = \langle V, \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle, \mathbf{f}, \epsilon, c_\alpha, c_r$ )

---

**Input:** a 3D truss  $\mathcal{T} = \langle V, \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$  with  $n$  vertices, which is a union of  $k$  convex edge-simple trusses  $\mathcal{T}_1, \dots, \mathcal{T}_k$ , a vector  $\mathbf{f} \in \mathbb{R}^{3n}$ , an error parameter  $\epsilon > 0$ .

Constants for aspect ratio threshold  $0 < c_\alpha < 1$ , and hollowing rate  $0 < c_r < 1$

**Output:** an approximate solution  $\mathbf{x}$  such that  $\|\mathbf{A}_{\mathcal{T}}\mathbf{x} - \mathbf{f}\|_2 \leq \epsilon \|\mathbf{f}\|_2$ .

1: **for** each  $i$  **do**

2:   Compute a bounding box  $B_i$  of  $\mathcal{T}_i$ .

3: **end for**

4: Let  $\mathcal{I} = \{1 \leq i \leq k : \alpha(\mathcal{T}_i) \leq n_i^{c_\alpha}\}$ .

5: **for** each  $i \in \mathcal{I}$  **do**

6:   Hollow out the interior vertices of  $\mathcal{T}_i$  with parameter  $r_i = n_i^{c_r}$  to form  $\mathcal{H}_i$ .

7: **end for**

8: Run nested dissection on the preconditioner (possibly with a specific set of separators).

9: Run preconditioned conjugate gradient with this preconditioner to solve the overall system.

10: **return** the solution  $\mathbf{x}$ .

---

**Bounding eigenvalues of an edge-simple and stiffly connected truss.** The key to analyze the running time of Algorithm 12 is to bound the relative condition number of the intermediate matrix and its preconditioner. By Theorem 2.5.3, this condition number determines the number of PCG iterations in line 9.

Since each vertex of a 3-D truss is incident to a constant number of edges with constant edge weights, it is not hard to check that its largest eigenvalue is some constant. The

following Lemma lower bounds its smallest non-zero eigenvalue.

**Lemma 5.1.1.** *Let  $\mathcal{T}$  be an edge-simple and stiffly-connected 3D truss. Let  $n$  be the number of vertices of  $\mathcal{T}$  and  $\Delta$  be the diameter. Let  $\mathbf{M}$  denote the associated stiffness matrix. Then,  $\lambda_{\min}(\mathbf{M}) = \Omega(n^{-1}\Delta^{-4})$  and  $\text{rank}(\mathbf{M}) = 3n - 6$ .*

The proof of Lemma 5.1.1 is quite involved, which is in Section 5.2 and 5.3. This proof is motivated by the Path Lemma in [44]. It bounds the smallest non-zero eigenvalue of a 2-D triangle path. In their proof, they center / normalize a vector in the eigen-space w.r.t. the first triangle of the path. Then they show that the minimum energy for stretching or squeezing each triangle to some position propagates along the triangle path, at a rate polynomial in the length of the path. However, this does not hold for 3-D tetrahedron surfaces. The propagated energy changes exponentially. We deal with this by trying different centerings, and show there always exists a centering desirable for our purpose.

We remark that we think the bound given in Lemma 5.1.1 is not tight, from the observations of Matlab experiments. It is not clear whether one can derive Cheeger type inequalities for truss stiffness matrices (even for the simplest 2-D case), similar to the one in [46].

**Proving the main result for small aspect ratio truss unions.** Motivated by the  $r$ -division introduced by Frederickson [77], we divide each small aspect ratio truss component into smaller *regions*. Given the convexity of each truss component, we can guarantee that the boundary of each small region is a well-shaped tetrahedron surface, and thus can be used as a good preconditioner of the Schur complement onto that boundary. The union of all these boundaries is called a *hollowing*. To create a hollowing, we fix two parameters: a bounding box  $B$  that determines the directions of each smaller chunks of the hollowing, and a size parameter  $r$  that controls the size of the smaller regions.

**Definition 5.1.2** ( $(B, r)$ -hollowing). *Given a convex edge-simple 3D truss  $\mathcal{T} = \langle V, \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$ , a bounding box  $B$  of  $\mathcal{T}$ , and a parameter  $r$ , a  $(B, r)$ -hollowing of  $\mathcal{T}$  is another edge-simple*



3D truss  $\mathcal{H} = \langle U, \{\mathbf{p}_i\}_{i \in U}, S, F, \gamma' \rangle$  such that,  $U \subseteq V$ ,  $S \subseteq T$ , and  $F$  is the subset of  $E$  that arises from edges in  $S$ , while  $\gamma'$  is just the restriction of  $\gamma$  to  $F$ . I.e. edges maintain the same stiffness factors as in  $\mathcal{T}$ . Also

1.  $\mathcal{H}$  contains  $O(nr^{-1/3})$  points.  $\mathcal{T} \setminus \mathcal{H}$  consists of  $O(nr^{-1})$  separate chunks, each of which has  $O(r)$  vertices and is incident to  $O(r^{2/3})$  vertices of  $\mathcal{H}$ .
2. for every plane  $P$  whose normal vector has angle  $\theta \in (0, \pi/2)$  with the longest direction of  $B$ , the number of tetrahedrons in  $\mathcal{H}$  intersected by  $P$  is

$$O\left(n^{2/3}\alpha^{-1/3}r^{-1/3}\cos^{-2}\theta\right),$$

where  $\alpha$  is the aspect ratio of  $\mathcal{T}$ .

3.  $\mathbf{A}_{\mathcal{H}} \preceq \text{SC}[\mathbf{A}_{\mathcal{T}}]_U \preceq O(r^2)\mathbf{A}_{\mathcal{H}}$ .

The next lemma describes the performance of algorithm HOLLOW, Algorithm 13 in Section 5.4, that we use to compute a  $(B, r)$ -hollowing of a convex edge-simple truss.

**Lemma 5.1.3.** *Given a convex edge-simple 3-D truss  $\mathcal{T} = \langle V, \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$  with  $n$  vertices, a bounding box  $B$  of  $\mathcal{T}$ , and a positive integer  $r$  such that the aspect ratio of  $\mathcal{T}$  is at most  $\sqrt{n/r}$ , the algorithm  $\text{HOLLOW}(\mathcal{T}, B, r)$  returns a  $(B, r)$ -hollowing  $\mathcal{H}$  of  $\mathcal{T}$ , and runs in time  $O(n)$ .*

We now prove our main result for small aspect ratio truss unions.

*Proof of Theorem 1.1.4.* Let  $\mathcal{T}$  be a edge-simple 3-D truss with  $n$  vertices, formed from a union of  $k$  convex edge-simple trusses, say  $\mathcal{T}_1, \dots, \mathcal{T}_k$ , each with aspect ratio at most  $O(1)$ . For each  $1 \leq i \leq k$ , let  $n_i$  be the number of vertices of  $\mathcal{T}_i$ , and define  $r_i \stackrel{\text{def}}{=} n_i^{1/2}$ . In Algorithm 12, for each  $\mathcal{T}_i$ , we compute a  $(B_i, r_i)$ -hollowing, where  $B_i$  is a bounding box

of  $\mathcal{T}_i$ . By Lemma 2.4.2 and 5.1.3, the total running time here is  $O(n)$ . In each  $(B_i, r_i)$ -hollowing region, we eliminate its interior vertices in total time

$$O\left(\sum_i n_i r_i^{-1} \cdot r_i^2\right) = O(n^{3/2}).$$

The Schur complement onto the boundaries has

$$O\left(\sum_i n_i r_i^{-1} \cdot (r_i^{2/3})^2\right) = O(n^{7/6})$$

nonzeros. We then run preconditioned conjugate gradient (PCG) to solve the linear system in the Schur complement by preconditioning it via the union of the  $(B_i, r_i)$ -hollowings, say  $\mathcal{T}'$ . Note by Jensen's inequality,  $\mathcal{T}'$  has size

$$O\left(\sum_i n_i r_i^{-1/3}\right) = O\left(\sum_i n_i^{5/6}\right) = O(k^{1/6} n^{5/6}).$$

Before running PCG, we compute a Cholesky factorization of  $\mathbf{A}_{\mathcal{T}'}$  by nested dissection. According to Theorem 2.5.1, the running time is  $O(k^{1/3} n^{5/3})$ , and the fill-in size is  $O(k^{2/9} n^{10/9})$ . By Definition 5.1.2, the condition number is  $O(\max_i r_i^2) = O(n)$ . According to Theorem 2.5.3, the number of PCG iterations is at most  $O(n^{1/2} \log(1/\epsilon))$  to output a solution up to accuracy  $\epsilon$ . In each PCG iteration, we do a matrix-vector multiplication with the Schur complement in time  $O(n^{7/6})$ , and solve a linear system in  $\mathbf{A}_{\mathcal{T}'}$  in time  $O(k^{2/9} n^{10/9})$ . Thus the total running time is  $O(k^{1/3} n^{5/3} \log(1/\epsilon))$ .  $\square$

**Proving the main result for all-aspect ratio truss unions.** We extend our result to cover the case when the union of convex edge-simple trusses also include trusses with arbitrarily large aspect ratios. The extension is based on an observation: large aspect ratio simplicial complexes imply the existence of small size separators. The following lemma specifies this observation, which we will prove in Section 5.5.2.

**Lemma 5.1.4.** *Given a convex edge-simple 3D truss  $\mathcal{T} = \langle V, \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$  with aspect ratio at least  $\alpha > 0$ , and its bounding box  $B$ . Let  $\mathbf{d} \in \mathbb{R}^3$  be a unit vector along the longest direction of  $B$ , and let  $\mathbf{g} \in \mathbb{R}^3$  be a unit vector with  $\mathbf{d} \cdot \mathbf{g} > 0$ . Then every plane orthogonal to  $\mathbf{g}$  intersects at most  $O(n^{2/3}\alpha^{-1/3}(\mathbf{d} \cdot \mathbf{g})^{-1})$  tetrahedrons.*

Combining this together with the existence of hollowings for small aspect ratio trusses, we get the following lemma. It speeds up nested dissection for a partial state of our Gaussian elimination, which is crucial to prove the main result for all-aspect ratio truss unions.

**Lemma 5.1.5 (Combining Separators).** *Given a edge-simple 3D truss  $\mathcal{T} = \langle V, \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$ , which is a union of  $k$  convex edge-simple trusses with up to  $n$  vertices in total. There exists a randomized algorithm which returns with high probability a truss  $\mathcal{T}'$  by selectively computing  $(B_i, r_i)$  hollowings of some of the pieces with parameter*

$$r_i \leq n_i^{1/3}$$

*so that a complete elimination of  $\mathcal{T}'$  has size  $O(n^{23/18}k^{44/9})$ , and takes time  $O(n^{11/6}k^{22/3})$  to compute.*

The algorithm that achieves Lemma 5.1.5 is Algorithm 15 CONVEXTRUSSUNIONND in Section 5.5. Given these lemmas, we can now sketch a proof of Theorem 1.1.5.

*Proof of Theorem 1.1.5.* We bound the running time of Algorithm 12 TRUSSSOLVER with the preconditioner and nested dissection constructed as per Lemma 5.1.5.

Since all the hollowings involve pieces with  $r_i \leq n_i^{1/3}$ , Definition 5.1.2 gives a bound of  $O(n^{2/3})$  on the condition number, and in turn a bound of  $O(n^{1/3} \log(1/\epsilon))$  on the number of PCG iterations via Theorem 2.5.3. Furthermore, similar to the proof of Theorem 1.1.4, the Schur complement of  $\mathcal{T}$  onto the elements of  $\mathcal{T}'$  has size  $O(n^{10/9})$ , and computing them by eliminating all interior vertices of our hollowings takes time  $O(n^{4/3})$ .

Thus, the total running time of Algorithm 12 is

$$O\left(n^{4/3} + n^{11/6}k^{22/3} + n^{1/3} \log(1/\epsilon) \cdot (n^{10/9} + n^{23/18}k^{44/9})\right) = O\left(n^{11/6}k^{22/3} \log(1/\epsilon)\right).$$

□

## 5.2 Bounding the Smallest Nonzero Eigenvalue of an edge-simple Truss

In this section, we prove Lemma 5.1.1, which lower bounds the smallest nonzero eigenvalue of an edge-simple 3D truss. We restate Lemma 5.1.1 in the following.

**Lemma 5.1.1.** *Let  $\mathcal{T}$  be an edge-simple and stiffly-connected 3D truss. Let  $n$  be the number of vertices of  $\mathcal{T}$  and  $\Delta$  be the diameter. Let  $\mathbf{M}$  denote the associated stiffness matrix. Then,  $\lambda_{\min}(\mathbf{M}) = \Omega(n^{-1}\Delta^{-4})$  and  $\text{rank}(\mathbf{M}) = 3n - 6$ .*

### 5.2.1 Main Ideas

The proof of Lemma 5.1.1 is an extension of the path support lemma by Daich and Spielman [44]. That proof relies on recentering a vector  $\mathbf{q}$ , which is a unit vector orthogonal to the null space, with respect to a single face by transforming it along the null space of  $\mathbf{M}$ .

Let  $\mathcal{T} = \langle V, \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$  be a 3D stiffly-connected truss over  $n$  vertices. The null space of the stiffness matrix of  $\mathcal{T}$  can be characterized as:

1.  $\mathbf{p}^x, \mathbf{p}^y, \mathbf{p}^z \in \mathbb{R}^{3n}$ : for each  $1 \leq i \leq n$ , the corresponding 3-dimensional vector  $\mathbf{p}_i^x$  ( $\mathbf{p}_i^y$  and  $\mathbf{p}_i^z$ ) has 1 for its  $x$ -coordinate ( $y$ -coordinate, and  $z$ -coordinate, respectively) and 0 for the other two coordinates.

2.  $\mathbf{p}^{\perp xy}, \mathbf{p}^{\perp xz}, \mathbf{p}^{\perp yz} \in \mathbb{R}^{3n}$ : fix an arbitrary index  $1 \leq c \leq n$ , for each  $1 \leq i \leq n$ :

$$\begin{aligned}\mathbf{p}_i^{\perp xy} &= \left[ -(\mathbf{p}_i - \mathbf{p}_c)_y, (\mathbf{p}_i - \mathbf{p}_c)_x, 0 \right]^\top, \\ \mathbf{p}_i^{\perp xz} &= \left[ -(\mathbf{p}_i - \mathbf{p}_c)_z, 0, (\mathbf{p}_i - \mathbf{p}_c)_x \right]^\top, \\ \mathbf{p}_i^{\perp yz} &= \left[ 0, -(\mathbf{p}_i - \mathbf{p}_c)_z, (\mathbf{p}_i - \mathbf{p}_c)_y \right]^\top.\end{aligned}$$

Also, as many of our arguments are symmetric across dimensions, we will use  $d, d_1$  and  $d_2$  to represent symmetric indexing over the dimensions, or pairs of dimensions respectively. Finally, as centering and exploring a simplicial complex from a particular triangle introduces an ordering on the tetrahedrons, faces, and edges, we will define our edges, triangles, and tetrahedrons as ordered tuples:

1. Edges:  $e = \langle e_1, e_2 \rangle$ ,
2. Triangles: we denote these as  $s = \langle s_1, s_2, s_3 \rangle$ . Here  $e(s)$  means the edge  $\langle s_1, s_2 \rangle$ .
3. Tetrahedrons: an ordered 4-tuples of pairwise adjacent points,  $t = \langle t_1, t_2, t_3, t_4 \rangle$ . Here  $s(t)$  means the (triangle) surface  $\langle t_1, t_2, t_3 \rangle$ , and  $e(t)$  means the edge  $e(s) = \langle t_1, t_2 \rangle$ .

With these notations in mind, we can center a vector  $\mathbf{q}$  w.r.t. a particular (oriented) triangle surface  $s$ . This can be viewed as an extension of the centering lemma in [44]. We prove the following Lemma in Section 5.2.2.

**Lemma 5.2.1.** *Given a stiffly-connected, edge-simple truss  $\mathcal{T} = \langle V, \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$ , let  $\mathbf{q}$  be a vector orthogonal to the null space of its stiffness matrix. For each oriented triangle  $s$ , there exists a (unique) vector  $\bar{\mathbf{q}}^{(s)}$  with scalar shift parameters  $c^{(s)\perp xy}, c^{(s)\perp xz}, c^{(s)\perp yz}$ :*

$$\bar{\mathbf{q}}^{(s)} = \mathbf{q} + \sum_{d_1 d_2} c^{(s)\perp d_1 d_2} \mathbf{p}^{\perp d_1 d_2} \quad (5.1)$$

satisfying:

1. the plane containing the points  $\bar{\mathbf{q}}_s^{(s)}$  is parallel to the plane containing  $\mathbf{p}_s$ .
2. The edge  $\bar{\mathbf{q}}_{e(s)}^{(s)}$  is parallel to the edge  $\mathbf{p}_{e(s)}$ .

Daitch and Spielman then showed that with any centering, the value of  $\mathbf{q}^T \mathbf{M} \mathbf{q}$  is lower bounded by the sum of a series of shifted values, or in simpler terms, the norm of  $\bar{\mathbf{q}}_i^{(s)} - \bar{\mathbf{q}}_j^{(s)}$  for some edge  $ij$ . However, our extension of this bound (which we will describe next in Lemma 5.2.4) to the 3-D case has an exponential dependency on the distance in tetrahedrons between  $ij$  and  $s$ . As a result, we first show the existence of a good centering, namely one where there exist an edge close to  $s$  whose endpoints are far apart. This notion of distance can be defined in terms of ‘hop count’ of tetrahedrons.

**Definition 5.2.2.** *The tetrahedron-distance between a pair of objects  $x$  and  $y$  in a simplicial complex is the shortest sequence of tetrahedrons*

$$t^{(0)}, t^{(1)}, \dots, t^{(d)}$$

*such that  $x \subseteq t^{(0)}$ ,  $y \subseteq t^{(d)}$ , and for all  $1 \leq i \leq d$ ,  $t^{(i-1)}$  and  $t^{(i)}$  share a triangle face.*

We remark that because all edges and angles are within some constant range, this combinatorial distance is within constant factors of the Euclidean distance of the associated points. However, we will not make use of this connection.

**Lemma 5.2.3.** *Given a stiffly-connected, edge-simple truss  $\mathcal{T}$  with  $n$  vertices and diameter  $\Delta$ , let  $\mathbf{q}$  be a unit vector orthogonal to the null space of the stiffness matrix of  $\mathcal{T}$ . There exists an oriented triangle  $s$  and a pair of points  $i, j$  within tetrahedron-distance  $O(1)$  of  $s$  satisfying:*

$$\left\| \bar{\mathbf{q}}_i^{(s)} - \bar{\mathbf{q}}_j^{(s)} \right\|_2^2 = \Omega \left( \frac{1}{\Delta^{4n}} \right).$$

We prove this lemma in Section 5.2.3.

We can then check, via an argument similar to [44], that such a centering and distance pair implies a large quadratic form. The following lemma will be proved in Section 5.3.

**Lemma 5.2.4.** *Given a stiffly-connected, edge-simple truss with stiffness matrix  $\mathbf{M}$ , an oriented triangle  $s$ , and a pair of points  $i, j$  within tetrahedron-distance  $h$  of  $s$ , we have*

$$\|\bar{\mathbf{q}}^{(s)}\|_{\mathbf{M}}^2 \geq 2^{-\Theta(h)} \left\| \bar{\mathbf{q}}_i^{(s)} - \bar{\mathbf{q}}_j^{(s)} \right\|_2^2.$$

*Proof of Lemma 5.1.1.* Consider an arbitrarily fixed unit vector  $\mathbf{q}$  that's orthogonal to the null space of  $\mathbf{M}$ . Let  $s$  be the centering given by Lemma 5.2.3, and let  $i, j$  be a pair of points within a constant tetrahedron-distance of  $s$ . Lemma 5.2.4 gives

$$\|\bar{\mathbf{q}}^{(s)}\|_{\mathbf{M}}^2 \geq \Omega(1) \left\| \bar{\mathbf{q}}_i^{(s)} - \bar{\mathbf{q}}_j^{(s)} \right\|_2^2 \geq \Omega\left(\frac{1}{\Delta^{4n}}\right). \quad (5.2)$$

On the other hand, by Equation (5.1),  $\|\bar{\mathbf{q}}^{(s)}\|_2 \geq \|\mathbf{q}\|_2 = 1$ . Thus,  $\lambda_{\min}(\mathbf{M}) \geq \Omega\left(\frac{1}{\Delta^{4n}}\right)$ .  $\square$

### 5.2.2 Centering a Vector (Proof of Lemma 5.2.1)

To prove Lemma 5.2.1, we define the following operation. Let  $P$  be any fixed plane in  $\mathbb{R}^3$  and let  $\mathbf{y} \in \mathbb{R}^3$ . We define  $\mathbf{y}^{\perp P}$  to be the vector obtained by first projecting  $\mathbf{y}$  onto plane  $P$  and then rotating the projected vector on the plane counterclockwise by  $\pi/2$ . The following claim shows that  $\mathbf{y}^{\perp P}$  can be written as a linear combination of  $\mathbf{y}^{\perp xy}, \mathbf{y}^{\perp yz}, \mathbf{y}^{\perp xz}$ .

**Claim 5.2.5** (Rotation matrix). *Let  $P$  be any fixed plane in  $\mathbb{R}^3$ , and let  $\mathbf{w}$  be its normal vector. Then,*

$$\forall \mathbf{y} \in \mathbb{R}^3, \quad \mathbf{y}^{\perp P} = \mathbf{w}_z \mathbf{y}^{\perp xy} - \mathbf{w}_y \mathbf{y}^{\perp xz} + \mathbf{w}_x \mathbf{y}^{\perp yz}.$$

*Proof.* Let

$$\mathbf{R} \stackrel{\text{def}}{=} \begin{pmatrix} 0 & -\mathbf{w}_z & \mathbf{w}_y \\ \mathbf{w}_z & 0 & -\mathbf{w}_x \\ -\mathbf{w}_y & \mathbf{w}_x & 0 \end{pmatrix}$$

be the rotation matrix which rotates a vector on the plane  $P$  counterclockwise by  $\pi/2$ . Then,

$$\mathbf{y}^{\perp_P} = \mathbf{R}(\mathbf{y} - (\mathbf{y}^\top \mathbf{w})\mathbf{w}) = \mathbf{R}(\mathbf{I} - \mathbf{w}^\top \mathbf{w})\mathbf{y}.$$

We can check that  $\mathbf{R}\mathbf{w} = \mathbf{0}$ . It implies that  $\mathbf{y}^{\perp_P} = \mathbf{R}\mathbf{y}$ .

On the other hand,

$$\mathbf{w}_z \mathbf{y}^{\perp_{xy}} - \mathbf{w}_y \mathbf{y}^{\perp_{xz}} + \mathbf{w}_x \mathbf{y}^{\perp_{yz}} = \begin{pmatrix} 0 & -\mathbf{w}_z & \mathbf{w}_y \\ \mathbf{w}_z & 0 & -\mathbf{w}_x \\ -\mathbf{w}_y & \mathbf{w}_x & 0 \end{pmatrix} \mathbf{y} = \mathbf{R}\mathbf{y}.$$

Thus, the claim holds. □

**Claim 5.2.6.** *Let  $\mathbf{h} \in \mathbb{R}^3$  be a nonzero vector. Then matrix*

$$\mathbf{H} \stackrel{\text{def}}{=} \begin{pmatrix} 0 & -\mathbf{h}_z & \mathbf{h}_y \\ \mathbf{h}_z & 0 & -\mathbf{h}_x \\ -\mathbf{h}_y & \mathbf{h}_x & 0 \\ \mathbf{h}_x & \mathbf{h}_y & \mathbf{h}_z \end{pmatrix}.$$

*has rank 3.*

*Proof.* The 2-by-2 bottom left submatrix has determinant  $-(\mathbf{h}_y^2 + \mathbf{h}_x^2)$ . If  $\mathbf{h}_y^2 + \mathbf{h}_x^2 = 0$ , then clearly  $\mathbf{H}$  has rank 3 and we have done; otherwise, the 3rd and the 4th rows are independent.

Now it suffices to show that the 2nd row is independent of the 3rd and the 4th rows.

Assume by contradiction, suppose

$$\mathbf{h}_z = -\alpha \mathbf{h}_y + \beta \mathbf{h}_x,$$

$$0 = \alpha \mathbf{h}_x + \beta \mathbf{h}_y,$$

$$-\mathbf{h}_x = \beta \mathbf{h}_z.$$



By solving the last two equations, we get  $\beta = -\mathbf{h}_x/\mathbf{h}_z$  and  $\alpha = \mathbf{h}_y/\mathbf{h}_z$ . Plugging these values into the 1st equation, we have  $\mathbf{h}_x^2 + \mathbf{h}_y^2 + \mathbf{h}_z^2 = 0$ , which contradicts that  $\mathbf{h} \neq \mathbf{0}$ .  $\square$

*Proof of Lemma 5.2.1.* Let  $\mathbf{w}$  be the normal vector of the plane containing  $s = \langle \mathbf{p}_{i_1}, \mathbf{p}_{i_2}, \mathbf{p}_{i_3} \rangle$ , that is,

$$\mathbf{w}^\top(\mathbf{p}_{i_2} - \mathbf{p}_{i_1}) = 0, \mathbf{w}^\top(\mathbf{p}_{i_3} - \mathbf{p}_{i_1}) = 0, \text{ and } \|\mathbf{w}\|_2 = 1.$$

We first show that there exist  $\alpha_{xy}, \alpha_{xz}, \alpha_{yz} \in \mathbb{R}$  such that the vector

$$\mathbf{z} \stackrel{\text{def}}{=} \mathbf{q} + \sum_{d1d2} \alpha_{d1d2} \mathbf{p}^{\perp_{d1d2}}$$

satisfies the first condition. It suffices to show that the following linear system has a solution for real numbers  $\alpha_{d1d2}$ 's:

$$\begin{aligned} \mathbf{w}^\top(\mathbf{q}_{i_2} - \mathbf{q}_{i_1} + \sum_{d1d2} \alpha_{d1d2}(\mathbf{p}_{i_2} - \mathbf{p}_{i_1})^{\perp_{d1d2}}) &= 0, \\ \mathbf{w}^\top(\mathbf{q}_{i_3} - \mathbf{q}_{i_1} + \sum_{d1d2} \alpha_{d1d2}(\mathbf{p}_{i_3} - \mathbf{p}_{i_1})^{\perp_{d1d2}}) &= 0. \end{aligned}$$

Rearrange it and write it in matrix form:

$$\begin{pmatrix} \mathbf{w}^\top(\mathbf{p}_{i_2} - \mathbf{p}_{i_1})^{\perp_{xy}} & \mathbf{w}^\top(\mathbf{p}_{i_2} - \mathbf{p}_{i_1})^{\perp_{yz}} & \mathbf{w}^\top(\mathbf{p}_{i_2} - \mathbf{p}_{i_1})^{\perp_{xz}} \\ \mathbf{w}^\top(\mathbf{p}_{i_3} - \mathbf{p}_{i_1})^{\perp_{xy}} & \mathbf{w}^\top(\mathbf{p}_{i_3} - \mathbf{p}_{i_1})^{\perp_{yz}} & \mathbf{w}^\top(\mathbf{p}_{i_3} - \mathbf{p}_{i_1})^{\perp_{xz}} \end{pmatrix} \begin{pmatrix} \alpha_{xy} \\ \alpha_{yz} \\ \alpha_{xz} \end{pmatrix} = \begin{pmatrix} -\mathbf{w}^\top(\mathbf{x}_{i_2} - \mathbf{x}_{i_1}) \\ -\mathbf{w}^\top(\mathbf{x}_{i_3} - \mathbf{x}_{i_1}) \end{pmatrix}.$$

It suffices to show that the coefficient matrix has rank 2.

Assume by contradiction, there is some  $k \in \mathbb{R}$  such that

$$\mathbf{w}^\top(\mathbf{p}_{i_2} - \mathbf{p}_{i_1})^{\perp_{d1d2}} = k \mathbf{w}^\top(\mathbf{p}_{i_3} - \mathbf{p}_{i_1})^{\perp_{d1d2}}, \quad \forall d1, d2$$

It equals to

$$\mathbf{w}^\top ((\mathbf{p}_{i_2} - \mathbf{p}_{i_1}) - k(\mathbf{p}_{i_3} - \mathbf{p}_{i_1}))^{\perp_{d1d2}} = 0, \quad \forall d1, d2$$

Let  $\mathbf{h} \stackrel{\text{def}}{=} (\mathbf{p}_{i_2} - \mathbf{p}_{i_1}) - k(\mathbf{p}_{i_3} - \mathbf{p}_{i_1})$ . Since vectors  $\mathbf{p}_{i_2} - \mathbf{p}_{i_1}, \mathbf{p}_{i_3} - \mathbf{p}_{i_1}$  are not parallel, we have  $\mathbf{h} \neq \mathbf{0}$ . Besides,  $\mathbf{h} \perp \mathbf{w}$ . Write the above equation in matrix form:

$$\begin{pmatrix} 0 & -\mathbf{h}_z & \mathbf{h}_y \\ \mathbf{h}_z & 0 & -\mathbf{h}_x \\ -\mathbf{h}_y & \mathbf{h}_x & 0 \\ \mathbf{h}_x & \mathbf{h}_y & \mathbf{h}_z \end{pmatrix} \mathbf{w} = \mathbf{0}.$$

By Claim 5.2.6, the coefficient matrix has rank 3, which implies that  $\mathbf{w} = \mathbf{0}$ . It contradicts that  $\|\mathbf{w}\|_2 = 1$ .

Then we show that there exist  $\beta_{xy}, \beta_{xz}, \beta_{yz} \in \mathbb{R}$  such that

$$\mathbf{q}^{(s)} = \mathbf{z} + \sum_{d1d2} \beta_{d1d2} \mathbf{p}^{\perp_{d1d2}}$$

satisfies both conditions.

Let  $P$  be the plane containing  $\mathbf{p}_{i_1}, \mathbf{p}_{i_2}, \mathbf{p}_{i_3}$ . Let  $\mathbf{g} \in \mathbb{R}^{3n}$  satisfy

$$\mathbf{g}_i = (\mathbf{p}_i - \mathbf{p}_1)^{\perp_P}, \quad \forall 1 \leq i \leq n$$

Then, there exists an appropriate multiplier  $\gamma \in \mathbb{R}$  such that the vector

$$(\mathbf{z}_{i_2} - \mathbf{z}_{i_1}) + \gamma(\mathbf{g}_{i_2} - \mathbf{g}_{i_1}) = (\mathbf{z}_{i_2} - \mathbf{z}_{i_1}) + \gamma(\mathbf{p}_{i_2} - \mathbf{p}_{i_1})^{\perp_P}$$

is parallel to  $\mathbf{p}_{i_2} - \mathbf{p}_{i_1}$ . Besides, since both  $\mathbf{z}$  and  $\mathbf{g}$  are parallel to the plane  $P$ , the vector  $\mathbf{q}^{(s)} = \mathbf{z} + \gamma\mathbf{g}$  is parallel to the plane  $P$ .

By Claim 5.2.5, there exists real numbers  $\beta_{xy}, \beta_{xz}, \beta_{yz}$  such that

$$\gamma \mathbf{g} = \sum_{d1d2} \beta_{d1d2} \mathbf{p}^{\perp_{d1d2}}.$$

Let

$$c^{\langle s \rangle \perp_{d1d2}} = \alpha_{d1d2} + \beta_{d1d2}, \quad \forall d1, d2$$

This completes the proof. □

### 5.2.3 Existence of Good Centering (Proof of Lemma 5.2.3)

In this section we prove Lemma 5.2.3.

The proof is by contradiction. We assume that for *every* centering at a triangle  $s$ , for every pair  $\langle i, j \rangle$  within tetrahedron-distance 3 of  $\langle s \rangle$  satisfies

$$\left\| \bar{\mathbf{q}}_i^{\langle s \rangle} - \bar{\mathbf{q}}_j^{\langle s \rangle} \right\|_2 \leq \epsilon, \tag{5.3}$$

where  $\epsilon \stackrel{\text{def}}{=} \sqrt{\frac{1}{\beta \Delta^4 n}}$  for a sufficiently large constant  $\beta$ .

Recall that in Lemma 5.2.1, for each centering at  $\langle s \rangle$ , we define 3 scalar coefficients

$$c^{\langle s \rangle \perp_{xy}}, c^{\langle s \rangle \perp_{xz}}, c^{\langle s \rangle \perp_{yz}} \in \mathbb{R}$$

for the 3 null space vectors in Equation (5.1). We will write the vector containing these 3 coefficients as  $\mathbf{c}^{\langle s \rangle}$ .

To prove Lemma 5.2.3, we need the following lemma. It says that, under the assumption in Equation (5.3), the difference between the coefficient vectors w.r.t. to two close centering triangles is small.

**Lemma 5.2.7.** *Assume that for every centering triangle  $s$  and every pair of points  $i, j$  within distance 3 of  $\langle s \rangle$  satisfies Equation (5.3). Then for every pair of centering triangles*

$s_1$  and  $s_2$  within constant tetrahedron distance to each other, we have

$$\|\mathbf{c}^{\langle s_1 \rangle} - \mathbf{c}^{\langle s_2 \rangle}\|_2 = O(\epsilon).$$

The next lemma implies that: for any two vertices such that each is centered w.r.t. a triangle close to itself, the difference between the two centered vertices is small.

**Lemma 5.2.8.** *Let  $u, w$  be two arbitrary vertices of  $\mathcal{T}$ . Let  $s_u$  ( $s_w$ ) be a triangle containing  $u$  (and  $w$ , respectively). Under the assumption in Equation (5.3), we have*

$$\|\bar{\mathbf{q}}_u^{\langle s_u \rangle} - \bar{\mathbf{q}}_w^{\langle s_w \rangle}\|_2 = O(\Delta^2 \epsilon).$$

*Proof.* Let  $u = v_1, v_2, \dots, v_f = w$  be a shortest path from  $u$  to  $w$ . Note  $f \leq \Delta$ . Let  $s_i$  be a triangle next to  $v_i$  for  $2 \leq i \leq f-1$ . The path from vertex  $u$ , centered at  $\langle s_u \rangle$ , to vertex  $w$ , centered at  $\langle s_w \rangle$ , can be expressed as the following:

$$\bar{\mathbf{q}}_u^{\langle s_u \rangle} - \bar{\mathbf{q}}_w^{\langle s_w \rangle} = \bar{\mathbf{q}}_u^{\langle s_u \rangle} - \bar{\mathbf{q}}_{v_2}^{\langle s_u \rangle} + \sum_{2 \leq i \leq f-1} (\bar{\mathbf{q}}_{v_i}^{\langle s_i \rangle} - \bar{\mathbf{q}}_{v_{i+1}}^{\langle s_{i+1} \rangle}) + \sum_{2 \leq i \leq f-1} (\bar{\mathbf{q}}_{v_i}^{\langle s_{i+1} \rangle} - \bar{\mathbf{q}}_{v_{i+1}}^{\langle s_{i+1} \rangle}).$$

Taking  $\ell_2$  norm on both sides and applying the triangle inequality, we can bound the norm of the LHS by the sum of  $\ell_2$  norm of each term in the RHS.

By Equation (5.1) and the triangle inequality,

$$\|\bar{\mathbf{q}}_{v_i}^{\langle s_i \rangle} - \bar{\mathbf{q}}_{v_i}^{\langle s_{i+1} \rangle}\|_2 \leq \sum_{d_1 d_2} |\mathbf{c}^{\langle s_i \rangle \perp d_1 d_2} - \mathbf{c}^{\langle s_{i+1} \rangle \perp d_1 d_2}| \|(\mathbf{p}_{v_i} - \mathbf{p}_c)^{\perp d_1 d_2}\|_2.$$

Apply Lemma 5.2.7:

$$\|\bar{\mathbf{q}}_{v_i}^{\langle s_i \rangle} - \bar{\mathbf{q}}_{v_i}^{\langle s_{i+1} \rangle}\|_2 = O(\Delta \epsilon). \quad (5.4)$$

Together with our assumption in Equation (5.3), we have

$$\|\bar{\mathbf{q}}_u^{(s_u)} - \bar{\mathbf{q}}_w^{(s_w)}\|_2 \leq O(\Delta^2 \epsilon).$$

□

Now we prove Lemma 5.2.3.

*Proof of Lemma 5.2.3.* For each vertex  $i$ , let  $s_i$  denote an arbitrary triangle next to vertex  $i$ .

We can write vector  $\mathbf{q}$  as

$$\mathbf{q} = \hat{\mathbf{q}} + \tilde{\mathbf{q}} - \sum_{d_1 d_2} c^{(s_1) \perp d_1 d_2} \mathbf{p}^{\perp d_1 d_2} + \mathbf{e}, \quad (5.5)$$

where

$$\hat{\mathbf{q}}_i = \bar{\mathbf{q}}_i^{(s_i)} - \bar{\mathbf{q}}_1^{(s_1)}, \tilde{\mathbf{q}}_i = \bar{\mathbf{q}}_i^{(s_1)} - \bar{\mathbf{q}}_i^{(s_i)}, \mathbf{e}_i = \bar{\mathbf{q}}_1^{(s_1)}.$$

Note that the last two terms of Equation (5.5) are in the null space of  $\mathbf{M}$ . Thus,

$$\|\hat{\mathbf{q}} + \tilde{\mathbf{q}}\|_2 \geq \|\mathbf{q}\|_2 = 1.$$

On the other hand, by the triangle inequality,

$$\|\hat{\mathbf{q}} + \tilde{\mathbf{q}}\|_2 \leq \|\hat{\mathbf{q}}\|_2 + \|\tilde{\mathbf{q}}\|_2 = \left( \sum_i \left\| \bar{\mathbf{q}}_i^{(s_i)} - \bar{\mathbf{q}}_1^{(s_1)} \right\|_2^2 \right)^{1/2} + \left( \sum_i \left\| \bar{\mathbf{q}}_i^{(s_1)} - \bar{\mathbf{q}}_i^{(s_i)} \right\|_2^2 \right)^{1/2}.$$

We apply Lemma 5.2.8 to the first term, and apply Equation (5.4) (which is true for any two close centering triangles)  $\Delta$  times for the second term:

$$\|\hat{\mathbf{q}} + \tilde{\mathbf{q}}\|_2 \leq \|\hat{\mathbf{q}}\|_2 + \|\tilde{\mathbf{q}}\|_2 = O(\sqrt{n} \Delta^2 \epsilon).$$

By our choice of  $\epsilon \leftarrow \frac{1}{\beta \sqrt{n} \Delta^2}$  for a sufficiently large constant  $\beta$ , we get a contradiction. □

It remains to prove Lemma 5.2.7. For it, we define a matrix w.r.t. a given vector  $\mathbf{v} \in \mathbb{R}^3$ :

$$\mathbf{Q}_v \stackrel{\text{def}}{=} (\mathbf{v}^{\perp yz}, \mathbf{v}^{\perp xz}, \mathbf{v}^{\perp xy}) = \begin{pmatrix} 0 & -v_z & v_y \\ v_z & 0 & -v_x \\ -v_y & v_x & 0 \end{pmatrix}. \quad (5.6)$$

We will use the following properties of  $\mathbf{Q}_v$ .

**Lemma 5.2.9.** *For the matrix  $\mathbf{Q}_v$  as defined in Equation (5.6), its singular values are 0, 1, 1, and its null space is multiples of the vector  $\mathbf{v}$ .*

*Proof.* For any vector  $\mathbf{p} \in \mathbb{R}^3$ , the cross product  $\mathbf{v} \times \mathbf{p} = \mathbf{Q}_v \mathbf{p}$ . That is,  $\mathbf{Q}_v \mathbf{p}$  is the vector obtained by:

1. First projecting  $\mathbf{p}$  onto the plane with normal vector  $\mathbf{v}$ , say  $P_{\perp v}$ , and
2. then rotating the projected vector on the plane  $P_{\perp v}$  by  $\frac{\pi}{2}$  counterclockwise.

From this description, we can infer that

$$\text{null}(\mathbf{Q}_v) = \mathbf{Span}(\mathbf{v})$$

and any vector orthogonal to  $\text{null}(\mathbf{Q}_v)$  is on this plane  $P_{\perp v}$ . Such vectors are not affected by the first step projection, and their lengths are not changed by the subsequent rotation. This means for such vectors  $\mathbf{u} \perp \text{null}(\mathbf{Q}_v)$  we have  $\mathbf{u}^\top \mathbf{Q}_v^\top \mathbf{Q}_v \mathbf{u} = \|\mathbf{u}\|_2^2$ . Thus, the singular values of  $\mathbf{Q}_v$  are 1,1,0.  $\square$

**Lemma 5.2.10.** *If  $\mathbf{v}$  and  $\mathbf{u}$  are vectors with length at least 1 such that the angle between them is  $\theta \in (0, \pi)$ , then the matrix*

$$\mathbf{Q}_v^\top \mathbf{Q}_v + \mathbf{Q}_u^\top \mathbf{Q}_u$$

*is full rank, and has minimum eigenvalue at least  $2 \sin^2 \frac{\theta}{2}$ .*

*Proof.* Let  $\mathbf{h} \in \mathbb{R}^3$  be a unit vector. Decompose  $\mathbf{h}$ :

$$\mathbf{h} = \alpha_v \mathbf{v} + \beta_v \hat{\mathbf{v}} = \alpha_u \mathbf{u} + \beta_u \hat{\mathbf{u}},$$

where  $\alpha_v, \alpha_u, \beta_v, \beta_u \in \mathbb{R}$ ,  $\hat{\mathbf{v}}$  is a unit vector orthogonal to  $\mathbf{v}$ , and  $\hat{\mathbf{u}}$  is a unit vector orthogonal to  $\mathbf{u}$ . Then,

$$\mathbf{h}^\top (\mathbf{Q}_v^\top \mathbf{Q}_v + \mathbf{Q}_u^\top \mathbf{Q}_u) \mathbf{h} = \beta_v^2 \|\mathbf{Q}_v \hat{\mathbf{v}}\|_2^2 + \beta_u^2 \|\mathbf{Q}_u \hat{\mathbf{u}}\|_2^2.$$

By Lemma 5.2.9,  $\|\mathbf{Q}_v \hat{\mathbf{v}}\|_2 = 1$  and  $\|\mathbf{Q}_u \hat{\mathbf{u}}\|_2 = 1$ . Thus,

$$\mathbf{h}^\top (\mathbf{Q}_v^\top \mathbf{Q}_v + \mathbf{Q}_u^\top \mathbf{Q}_u) \mathbf{h} = \beta_v^2 + \beta_u^2 = 2 - (\alpha_v^2 + \alpha_u^2).$$

The second equality is due to that  $\mathbf{h}$  is a unit vector.

$$\alpha_v^2 + \alpha_u^2 = \|\mathbf{h}^\top \mathbf{v}\|_2^2 + \|\mathbf{h}^\top \mathbf{u}\|_2^2 \leq 2 \cos^2 \frac{\theta}{2}.$$

Therefore,

$$\lambda_{\min}(\mathbf{Q}_v^\top \mathbf{Q}_v + \mathbf{Q}_u^\top \mathbf{Q}_u) \geq 2 \sin^2 \frac{\theta}{2}.$$

$\theta < \pi$  implies that  $\lambda_{\min}(\mathbf{Q}_v^\top \mathbf{Q}_v + \mathbf{Q}_u^\top \mathbf{Q}_u) > 0$ , and thus the matrix  $\mathbf{Q}_v^\top \mathbf{Q}_v + \mathbf{Q}_u^\top \mathbf{Q}_u$  has full rank. □

Equipped with the above lemmas, we prove Lemma 5.2.7.

*Proof of Lemma 5.2.7.* Let  $s_1$  and  $s_2$  be two triangle centerings for which there exist two edges belong to same tetrahedron, say  $(i, j), (j, k)$ , such that vertices  $i, j, k$  are all within tetrahedron-distance constant  $h$  of both  $s_1$  and  $s_2$ .

Recall  $\overline{\mathbf{q}}^{\langle s_1 \rangle}$  is defined in Equation (5.1). Subtracting

$$\overline{\mathbf{q}}_i^{\langle s_1 \rangle} = \mathbf{q}_i + \sum_{d_1 d_2} c^{\langle s_1 \rangle \perp d_1 d_2} (\mathbf{p}_i^{\perp d_1 d_2} - \mathbf{p}_c^{\perp d_1 d_2})$$

from the corresponding equation for  $j$  gives:

$$\overline{\mathbf{q}}_i^{\langle s_1 \rangle} - \overline{\mathbf{q}}_j^{\langle s_1 \rangle} = \mathbf{q}_i - \mathbf{q}_j + \sum_{d_1 d_2} c^{\langle s_1 \rangle \perp d_1 d_2} (\mathbf{p}_i^{\perp d_1 d_2} - \mathbf{p}_j^{\perp d_1 d_2}).$$

Plugging in the definition of  $\mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)}$ :

$$\overline{\mathbf{q}}_i^{\langle s_1 \rangle} - \overline{\mathbf{q}}_j^{\langle s_1 \rangle} = \mathbf{q}_i - \mathbf{q}_j + \mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)} \mathbf{c}^{\langle s_1 \rangle}.$$

Subtracting this equation from centering triangle  $s_2$  in turn cancels the  $\mathbf{q}_i - \mathbf{q}_j$  term on the RHS, giving:

$$\mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)} (\mathbf{c}^{\langle s_1 \rangle} - \mathbf{c}^{\langle s_2 \rangle}) = \left( \overline{\mathbf{q}}_i^{\langle s_1 \rangle} - \overline{\mathbf{q}}_j^{\langle s_1 \rangle} \right) - \left( \overline{\mathbf{q}}_i^{\langle s_2 \rangle} - \overline{\mathbf{q}}_j^{\langle s_2 \rangle} \right).$$

Together with the corresponding equation for  $j, k$ ,

$$\begin{pmatrix} \mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)} \\ \mathbf{Q}_{(\mathbf{p}_j - \mathbf{p}_k)} \end{pmatrix} (\mathbf{c}^{\langle s_1 \rangle} - \mathbf{c}^{\langle s_2 \rangle}) = \begin{pmatrix} \left( \overline{\mathbf{q}}_i^{\langle s_1 \rangle} - \overline{\mathbf{q}}_j^{\langle s_1 \rangle} \right) - \left( \overline{\mathbf{q}}_i^{\langle s_2 \rangle} - \overline{\mathbf{q}}_j^{\langle s_2 \rangle} \right) \\ \left( \overline{\mathbf{q}}_j^{\langle s_1 \rangle} - \overline{\mathbf{q}}_k^{\langle s_1 \rangle} \right) - \left( \overline{\mathbf{q}}_j^{\langle s_2 \rangle} - \overline{\mathbf{q}}_k^{\langle s_2 \rangle} \right) \end{pmatrix}.$$

Multiplying  $\begin{pmatrix} \mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)}^\top & \mathbf{Q}_{(\mathbf{p}_j - \mathbf{p}_k)}^\top \end{pmatrix}$  on both sides gives:

$$\begin{aligned} & \left( \mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)}^\top \mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)} + \mathbf{Q}_{(\mathbf{p}_j - \mathbf{p}_k)}^\top \mathbf{Q}_{(\mathbf{p}_j - \mathbf{p}_k)} \right) (\mathbf{c}^{\langle s_1 \rangle} - \mathbf{c}^{\langle s_2 \rangle}) \\ &= \begin{pmatrix} \mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)}^\top & \mathbf{Q}_{(\mathbf{p}_j - \mathbf{p}_k)}^\top \end{pmatrix} \begin{pmatrix} \left( \overline{\mathbf{q}}_i^{\langle s_1 \rangle} - \overline{\mathbf{q}}_j^{\langle s_1 \rangle} \right) - \left( \overline{\mathbf{q}}_i^{\langle s_2 \rangle} - \overline{\mathbf{q}}_j^{\langle s_2 \rangle} \right) \\ \left( \overline{\mathbf{q}}_j^{\langle s_1 \rangle} - \overline{\mathbf{q}}_k^{\langle s_1 \rangle} \right) - \left( \overline{\mathbf{q}}_j^{\langle s_2 \rangle} - \overline{\mathbf{q}}_k^{\langle s_2 \rangle} \right) \end{pmatrix}. \end{aligned}$$



Solving the above linear equations and taking norm on both sides:

$$\begin{aligned} \|\mathbf{c}^{\langle s_1 \rangle} - \mathbf{c}^{\langle s_2 \rangle}\|_2 &\leq \left\| \left( \mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)}^\top \mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)} + \mathbf{Q}_{(\mathbf{p}_j - \mathbf{p}_k)}^\top \mathbf{Q}_{(\mathbf{p}_j - \mathbf{p}_k)} \right)^{-1} \right\|_2 \\ &\quad \cdot \left\| \begin{pmatrix} \mathbf{Q}_{(\mathbf{p}_i - \mathbf{p}_j)}^\top & \mathbf{Q}_{(\mathbf{p}_j - \mathbf{p}_k)}^\top \end{pmatrix} \right\|_2 \left\| \begin{pmatrix} \left( \bar{\mathbf{q}}_i^{\langle s_1 \rangle} - \bar{\mathbf{q}}_j^{\langle s_1 \rangle} \right) - \left( \bar{\mathbf{q}}_i^{\langle s_2 \rangle} - \bar{\mathbf{q}}_j^{\langle s_2 \rangle} \right) \\ \left( \bar{\mathbf{q}}_j^{\langle s_1 \rangle} - \bar{\mathbf{q}}_k^{\langle s_1 \rangle} \right) - \left( \bar{\mathbf{q}}_j^{\langle s_2 \rangle} - \bar{\mathbf{q}}_k^{\langle s_2 \rangle} \right) \end{pmatrix} \right\|_2. \end{aligned}$$

Applying Lemma 5.2.9 and 5.2.10 on the first two terms, and applying the triangle inequality and the assumption in Equation (5.3) give:

$$\|\mathbf{c}^{\langle s_1 \rangle} - \mathbf{c}^{\langle s_2 \rangle}\|_2 \leq O(\epsilon).$$

□

### 5.3 Proof of Path Lemma

We now prove Lemma 5.2.4, which lower bounds the quadratic form of an edge-simple and stiffly-connected truss stiffness matrix by the distance between two centered points. As we now only deal with a single centering in this section, we will drop this superscription for simplicity and relabel all indices w.r.t. this centering. Our relabeling is similar to that in [44].

#### 5.3.1 Relabeling Tetrahedrons and Vertices

Fix an arbitrary tetrahedron, say  $t_1$ , and we center the vector  $\mathbf{q}$  w.r.t. one of the triangle faces of  $t_1$  as in Lemma 5.2.1. Use  $t_1$  as root, we run breadth-first-search (BFS) in the rigidity graph (refer to Definition 2.4.1) of  $\mathcal{T}$ , and relabel the tetrahedrons of  $\mathcal{T}$  according to this BFS ordering. For example, the neighbor tetrahedrons of  $t_1$  are labeled as  $t_2, t_3, \dots$ . Let  $T_{\text{BFS}}$  be the corresponding BFS tree in which each node represents a tetrahedron in  $\mathcal{T}$ .

Based on  $T_{\text{BFS}}$ , we relabel the vertices of  $\mathcal{T}$  as follows. We label the vertices of  $t_1$  by

$-2, -1, 0, 1$  in an arbitrary order. Each child of  $t_1$  (and subsequent recursions) shares a triangle face with their respective parent tetrahedron and thus only requires us to label one vertex per child tetrahedron, which can be labeled according to the BFS ordering.

For each newly labeled vertex  $j$ , we use  $t_j$  to denote the tetrahedron encompassing the  $t_j$  and its parental face. Besides, we use a 3-dimensional vector  $\sigma_j$  to consist of the indexes of the parental face sorted in ascending order, and is assumed that indexing this vector is implicitly modulo 3. Then,  $t_j = \{j, \sigma_j(1), \sigma_j(2), \sigma_j(3)\}$ . For completeness, we define  $\sigma_1 \stackrel{\text{def}}{=} \{-2, -1, 0\}$ . See Figure 5.1 for an example.

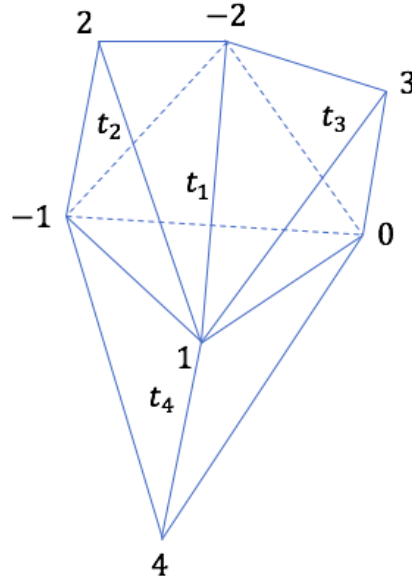


Figure 5.1: An example of the relabeling of tetrahedrons and vertices:  $t_1 = \{-2, -1, 0, 1\}$ ,  $t_2 = \{-2, -1, 1, 2\}$ ,  $t_3 = \{-2, 0, 1, 3\}$ ,  $t_4 = \{-1, 0, 1, 4\}$  and  $\sigma_1 = \{-2, -1, 0\}$ ,  $\sigma_2 = \{-2, -1, 1\}$ ,  $\sigma_3 = \{-2, 1, 0\}$ ,  $\sigma_4 = \{0, -1, 1\}$ .

### 5.3.2 Distance between Local Minimizers and the Centered Vectors

Let  $\mathbf{q} \in \mathbb{R}^{3n}$  be a unit vector which is orthogonal to  $\text{Span}(\mathbf{p}^x, \mathbf{p}^y, \mathbf{p}^z, \mathbf{p}^{\perp xy}, \mathbf{p}^{\perp yz}, \mathbf{p}^{\perp xz})$

defined at the beginning of Section 5.2.1. By Lemma 5.2.1, there exist scalars  $c^{\langle t_1 \rangle \perp xy}, c^{\langle t_1 \rangle \perp xz}, c^{\langle t_1 \rangle \perp yz} \in$

$\mathbb{R}$  such that the vector

$$\overline{\mathbf{q}}^{\langle t_1 \rangle} \stackrel{\text{def}}{=} \mathbf{q} + \sum_{d_1 d_2} c^{\langle t_1 \rangle \perp d_1 d_2} \mathbf{p}^{\perp d_1 d_2}$$

satisfies:

1. the plane containing  $\overline{\mathbf{q}}_{-2}^{(t_1)}, \overline{\mathbf{q}}_{-1}^{(t_1)}, \overline{\mathbf{q}}_0^{(t_1)}$  is parallel to the plane  $\sigma_{t_1}$ , and
2.  $\overline{\mathbf{q}}_{-2}^{(t_1)} - \overline{\mathbf{q}}_{-1}^{(t_1)}$  is parallel to  $(\sigma_{t_1}(1) - \sigma_{t_1}(2))$ .

We drop the superscription  $\langle t_1 \rangle$  when the context is clear.

For each  $0 \leq i \leq n-3$ , we define a 3-dimensional vector  $\mathbf{y}_i$ .  $\mathbf{y}_0$  is a vector on the plane of  $\sigma_{t_1} = \{\mathbf{p}_{-2}, \mathbf{p}_{-1}, \mathbf{p}_0\}$  and minimizes the energy / quadratic form, suppose both  $\overline{\mathbf{q}}_{-2}^{(t_1)}, \overline{\mathbf{q}}_{-1}^{(t_1)}$  are fixed. That is,

$$\begin{aligned} (\mathbf{p}_0 - \mathbf{p}_j)^\top (\mathbf{y}_0 - \overline{\mathbf{q}}_j) &= 0, \quad \forall j \in \{-2, -1\} \\ \mathbf{w}^\top (\mathbf{y}_0 - \overline{\mathbf{q}}_{-2}) &= 0 \end{aligned} \quad (5.7)$$

For each  $1 \leq i \leq n-3$ ,  $\mathbf{y}_i$  is a vector which minimizes the energy / quadratic form, suppose all the three points of  $\sigma_j$  are fixed. That is,

$$(\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)})^\top (\mathbf{y}_i - \overline{\mathbf{q}}_{\sigma_i(j)}) = 0, \quad \forall j \in \{1, 2, 3\} \quad (5.8)$$

We then define the distance between each local minimizer  $\mathbf{y}_i$  and the centered vector  $\overline{\mathbf{q}}_i$ . Specifically,

$$\mathbf{d}_i \stackrel{\text{def}}{=} \begin{cases} \overline{\mathbf{q}}_{-1} - \overline{\mathbf{q}}_{-2}, & i = -1 \\ \overline{\mathbf{q}}_i - \mathbf{y}_i, & 0 \leq i \leq n-3 \end{cases} \quad (5.9)$$

This definition intermediately gives that  $\forall -1 \leq i \leq n-3, 1 \leq j \leq \min\{i+2, 3\}$ ,

$$\mathbf{d}_i^\top (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)}) = (\overline{\mathbf{q}}_i - \mathbf{y}_i)^\top (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)}) = (\overline{\mathbf{q}}_i - \overline{\mathbf{q}}_{\sigma_i(j)})^\top (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)}). \quad (5.10)$$

Note that equation (5.10) describes the net stress on an edge.

We can explicitly express these  $\mathbf{d}_i$ 's by solving Equations (5.7) and (5.8).

Let

$$\mathbf{F}_0 \stackrel{\text{def}}{=} \begin{pmatrix} (\mathbf{p}_0 - \mathbf{p}_{-2})^\top \\ (\mathbf{p}_0 - \mathbf{p}_{-1})^\top \\ \mathbf{w}^\top \end{pmatrix}.$$

Let  $\mathbf{h}_{-2}, \mathbf{h}_{-1}, \mathbf{h}_w \in \mathbb{R}^3$  such that for each  $1 \leq k \leq 3$ ,

$$\begin{aligned} \mathbf{h}_{-2}(k) &\stackrel{\text{def}}{=} \det \begin{pmatrix} (\mathbf{p}_0 - \mathbf{p}_{-1})^{(k+1)_3} & (\mathbf{p}_0 - \mathbf{p}_{-1})^{(k+2)_3} \\ \mathbf{w}^{(k+1)_3} & \mathbf{w}^{(k+2)_3} \end{pmatrix}, \\ \mathbf{h}_{-1}(k) &\stackrel{\text{def}}{=} \det \begin{pmatrix} (\mathbf{p}_0 - \mathbf{p}_{-2})^{(k+1)_3} & (\mathbf{p}_0 - \mathbf{p}_{-2})^{(k+2)_3} \\ \mathbf{w}^{(k+1)_3} & \mathbf{w}^{(k+2)_3} \end{pmatrix}, \\ \mathbf{h}_w(k) &\stackrel{\text{def}}{=} \det \begin{pmatrix} (\mathbf{p}_0 - \mathbf{p}_{-1})^{(k+1)_3} & (\mathbf{p}_0 - \mathbf{p}_{-1})^{(k+2)_3} \\ (\mathbf{p}_0 - \mathbf{p}_{-2})^{(k+1)_3} & (\mathbf{p}_0 - \mathbf{p}_{-2})^{(k+2)_3} \end{pmatrix}. \end{aligned}$$

where  $(k+1)_3$  and  $(k+2)_3$  denotes accessing specific dimensions of vectors  $\in \mathbb{R}^3$  modulo

3. Then we define

$$\mathbf{H}_{0,2} \stackrel{\text{def}}{=} \mathbf{h}_{-2}(\mathbf{p}_0 - \mathbf{p}_{-2})^\top, \quad \mathbf{H}_{0,1} \stackrel{\text{def}}{=} \mathbf{h}_{-1}(\mathbf{p}_0 - \mathbf{p}_{-1})^\top, \quad \text{and } \mathbf{H}_w \stackrel{\text{def}}{=} \mathbf{h}_w \mathbf{w}^\top.$$

Similarly, for  $1 \leq i \leq n-3$ , let

$$\mathbf{F}_i \stackrel{\text{def}}{=} \begin{pmatrix} (\mathbf{p}_i - \mathbf{p}_{\sigma_i(1)})^\top \\ (\mathbf{p}_i - \mathbf{p}_{\sigma_i(2)})^\top \\ (\mathbf{p}_i - \mathbf{p}_{\sigma_i(3)})^\top \end{pmatrix},$$

and

$$\mathbf{H}_{i,j} \stackrel{\text{def}}{=} \mathbf{h}_{i,j}(\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)})^\top, \quad \forall 1 \leq j \leq 3, \quad (5.11)$$

where  $\mathbf{h}_{i,j} \in \mathbb{R}^3$  with

$$\mathbf{h}_{i,j}(k) \stackrel{\text{def}}{=} \det \begin{pmatrix} (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+1)_3)})^{(k+1)_3} & (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+1)_3)})^{(k+2)_3} \\ (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+2)_3)})^{(k+1)_3} & (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+2)_3)})^{(k+2)_3} \end{pmatrix}.$$

We can check that the following  $\mathbf{y}_i$ 's satisfy Equations (5.7) and (5.8):

$$\mathbf{y}_i = \begin{cases} \frac{1}{\det(\mathbf{F}_0)} ((\mathbf{H}_{0,2} + \mathbf{H}_w)\bar{\mathbf{q}}_{-2} + \mathbf{H}_{0,1}\bar{\mathbf{q}}_{-1}), & i = 0 \\ \frac{1}{\det(\mathbf{F}_i)} \sum_{1 \leq j \leq 3} \mathbf{H}_{i,j} \bar{\mathbf{q}}_{\sigma_i(j)}, & \forall 1 \leq i \leq n \end{cases}$$

**Claim 5.3.1.**

$$\det(\mathbf{F}_i) = (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)})^\top \mathbf{h}_{i,j}. \quad (5.12)$$

*Proof.* Let

$$\mathbf{Q} = \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix}$$

be a  $3 \times 3$  matrix. Then, the determinant of  $\mathbf{Q}$  is

$$\det(\mathbf{Q}) = a \cdot \det \left( \begin{pmatrix} e & f \\ h & i \end{pmatrix} \right) - b \cdot \det \left( \begin{pmatrix} d & f \\ g & i \end{pmatrix} \right) + c \cdot \det \left( \begin{pmatrix} d & e \\ g & h \end{pmatrix} \right).$$

Applying the above rule to  $\mathbf{F}_i$  gives Equation (5.12).  $\square$

**Claim 5.3.2.**

$$\det(\mathbf{F}_i) \mathbf{I} = \sum_{1 \leq j \leq 3} \mathbf{H}_{i,j}, \quad \forall i \geq 0$$

*Proof.* We first show that the diagonals of  $\sum_{j \in [3]} \mathbf{H}_{i,j}$  are equal to  $\det(\mathbf{F}_i)$ .

$$\begin{aligned}
& \sum_{j \in [3]} \mathbf{H}_{i,j}(1, 1) \\
&= \sum_{j \in [3]} (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)})(1) \cdot \det \begin{pmatrix} (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+1)_3)})(2) & (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+1)_3)})(3) \\ (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+2)_3)})(2) & (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+2)_3)})(3) \end{pmatrix} \\
&= \det \begin{pmatrix} (\mathbf{p}_i - \mathbf{p}_{\sigma_i(1)})^\top \\ (\mathbf{p}_i - \mathbf{p}_{\sigma_i(2)})^\top \\ (\mathbf{p}_i - \mathbf{p}_{\sigma_i(3)})^\top \end{pmatrix} \\
&= \det(\mathbf{F}_i).
\end{aligned}$$

Similarly, we have

$$\sum_{j \in [3]} \mathbf{H}_{i,j}(2, 2) = \sum_{j \in [3]} \mathbf{H}_{i,j}(3, 3) = \det(\mathbf{F}_i).$$

Then, we show that the off-diagonals of  $\sum_{j \in [3]} \mathbf{H}_{i,j}$  are 0.

$$\begin{aligned}
& \sum_{j \in [3]} \mathbf{H}_{i,j}(1, 2) \\
&= \sum_{j \in [3]} (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+1)_3)})(2) \cdot \det \begin{pmatrix} (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+1)_3)})(2) & (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+1)_3)})(3) \\ (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+2)_3)})(2) & (\mathbf{p}_i - \mathbf{p}_{\sigma_i((j+2)_3)})(3) \end{pmatrix} \\
&= \det \begin{pmatrix} (\mathbf{p}_i - \mathbf{p}_{\sigma_i(1)})(2) & (\mathbf{p}_i - \mathbf{p}_{\sigma_i(1)})(2) & (\mathbf{p}_i - \mathbf{p}_{\sigma_i(1)})(3) \\ (\mathbf{p}_i - \mathbf{p}_{\sigma_i(2)})(2) & (\mathbf{p}_i - \mathbf{p}_{\sigma_i(2)})(2) & (\mathbf{p}_i - \mathbf{p}_{\sigma_i(2)})(3) \\ (\mathbf{p}_i - \mathbf{p}_{\sigma_i(3)})(2) & (\mathbf{p}_i - \mathbf{p}_{\sigma_i(3)})(2) & (\mathbf{p}_i - \mathbf{p}_{\sigma_i(3)})(3) \end{pmatrix} \\
&= 0.
\end{aligned}$$

The last equation is due to the 3 columns of the matrix are linearly dependent. Similarly,

$$\sum_{j \in [3]} \mathbf{H}_{i,j}(s, t) = 0, \quad \forall 1 \leq s \neq t \leq 3.$$

This completes the proof. □

Plugging the above equations into the definition of  $\mathbf{d}_i$ 's gives:

**Claim 5.3.3.** *For each  $-1 \leq i \leq n-3$ ,*

$$\mathbf{d}_i = \bar{\mathbf{q}}_i - \bar{\mathbf{q}}_{\sigma_i(j)} + \frac{1}{\det(\mathbf{F}_i)} \sum_{\substack{1 \leq j' \leq \min\{i+2, 3\} \\ j' \neq j}} \mathbf{H}_{i,j'} (\bar{\mathbf{q}}_{\sigma_i(j)} - \bar{\mathbf{q}}_{\sigma_i(j')}).$$

### 5.3.3 Bounding the Norm of $\bar{\mathbf{q}}_i - \bar{\mathbf{q}}_{\sigma_i(j)}$ in terms of $\mathbf{d}_i$ 's

Rearranging the above equation gives that for each  $-1 \leq i \leq n-3, 1 \leq j \leq \min\{i+2, 3\}$ ,

$$\bar{\mathbf{q}}_i - \bar{\mathbf{q}}_{\sigma_i(j)} = \mathbf{d}_i - \frac{1}{\det(\mathbf{F}_i)} \sum_{\substack{1 \leq j' \leq \min\{i+2, 3\} \\ j' \neq j}} \mathbf{H}_{i,j'} (\bar{\mathbf{q}}_{\sigma_i(j)} - \bar{\mathbf{q}}_{\sigma_i(j')}). \quad (5.13)$$

Our goal is to express each  $\bar{\mathbf{q}}_i - \bar{\mathbf{q}}_{\sigma_i(j)}$  as a function of  $\mathbf{d}_i, \mathbf{d}_{i-1}, \dots, \mathbf{d}_{-1}$ . If each term  $\bar{\mathbf{q}}_{\sigma_i(j)} - \bar{\mathbf{q}}_{\sigma_i(j')}$  in the right hand side of Equation (5.13) satisfies

$$\{\sigma_i(j), \sigma_i(j')\} = \{i_1, \sigma_{i_1}(j_1)\}$$

for some  $i_1 < i$  and  $1 \leq j_1 \leq \min\{i_1+2, 3\}$ , then we can substitute it by Equation (5.13) with the left hand side being  $\bar{\mathbf{q}}_{i_1} - \bar{\mathbf{q}}_{\sigma_{i_1}(j_1)}$ . The substitution terminates when the right hand side term becomes  $\pm (\bar{\mathbf{q}}_{-1} - \bar{\mathbf{q}}_{-2}) = \pm \mathbf{d}_{-1}$ .

**Claim 5.3.4.** *Let  $\bar{\mathbf{q}}_{\sigma_{i_1}(j_1)} - \bar{\mathbf{q}}_{\sigma_{i_1}(j_2)}$  ( $0 \leq i_1 \leq n-3, 1 \leq j_1 \neq j_2 \leq \min\{i_1+2, 3\}$ ) be a term appearing in the right hand side of Equation (5.13). There exist  $-1 \leq i_2 < i_1, 1 \leq$*

$j_3 \leq \min\{i_2 + 2, 3\}$  satisfying

$$\{\sigma_{i_1}(j_1), \sigma_{i_1}(j_2)\} = \{i_2, \sigma_{i_2}(j_3)\}.$$

*Proof.* Without the loss of generality, assume  $\sigma_{i_1}(j_1) > \sigma_{i_1}(j_2)$ . If  $i_1 = 0$ , then  $\{\sigma_{i_1}(j_1), \sigma_{i_1}(j_2)\} = \{-1, \sigma_{-1}(1)\}$ ; if  $i_1 = 1$ , then  $\{\sigma_{i_1}(j_1), \sigma_{i_1}(j_2)\} \in \{\{0, \sigma_0(1)\}, \{0, \sigma_0(2)\}, \{-1, \sigma_0(1)\}\}$ .

The remaining proof focuses on  $i_1 \geq 2$ .

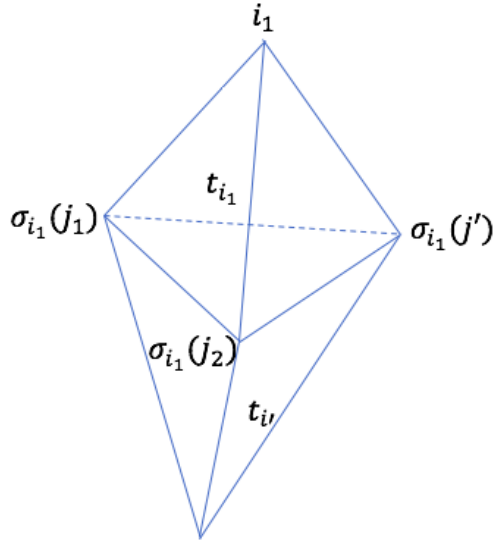


Figure 5.2: Tetrahedrons  $t_{i_1}$  and  $t_{i'}$

Let  $\sigma_{i_1}(j')$  be the other vertex in tetrahedron  $t_{i_1}$ . Let  $t_{i'}$  be the parent of  $t_{i_1}$  in the BFS tree  $T_{\text{BFS}}$ .  $i' < i_1$ , and the two tetrahedrons  $t_{i_1}$  and  $t_{i'}$  share a triangle face containing vertices  $\sigma_{i_1}(j_1), \sigma_{i_1}(j_2), \sigma_{i_1}(j')$ . See Figure 5.2. If  $i' = 1$ , then

$$\{\sigma_{i_1}(j_1), \sigma_{i_1}(j_2)\} \in \{\{i_2, \sigma_{i_2}(j_3)\} : -1 \leq i_2 \leq 1, 1 \leq j_3 \leq \min\{i_2 + 2, 3\}\}.$$

Otherwise  $i' \geq 2$ , we prove the statement by case analysis.

**Case 1.**  $\sigma_{i_1}(j_1) = \max_{1 \leq k \leq 3} \{\sigma_{i_1}(k)\}$ . By our labeling rules, vertex  $i'$  is the one with the



maximum index among all the vertices of  $t_{i'}$  and it is contained in the triangle shared by  $t_{i_1}$  and  $t_{i'}$ . Thus,  $\sigma_{i_1}(j_1) = i'$  and

$$\{\sigma_{i_1}(j_1), \sigma_{i_1}(j_2)\} = \{i', \sigma_{i'}(j_3)\}$$

for some  $1 \leq j_3 \leq 3$ .

**Case 2.**  $\sigma_{i_1}(j') = \max_{1 \leq k \leq 3} \{\sigma_{i_1}(k)\}$ . Then,  $\sigma_{i_1}(j') = i'$  and  $\sigma_{i_1}(j_1), \sigma_{i_1}(j_2) \in \{\sigma_{i'}(k) : 1 \leq k \leq 3\}$ . Note  $i' < i$ . By induction on the tetrahedron index  $i'$ , we can see that there exist  $i_2, j_3$  satisfying  $\{\sigma_{i_1}(j_1), \sigma_{i_1}(j_2)\} = \{i_2, \sigma_{i_2}(j_3)\}$ .  $\square$

We use a recursion tree  $T_{\text{rec}}^{(i,j)}$  to express the process of recursively substituting  $\bar{q}_{i'} - \bar{q}_{\sigma_{i'}(j')}$  via Equation (5.13). The root of  $T_{\text{rec}}^{(i,j)}$  is  $\bar{q}_i - \bar{q}_{\sigma_i(j)}$ , and each node  $u$  of  $T_{\text{rec}}^{(i,j)}$  represents a term  $\bar{q}_{i_u} - \bar{q}_{\sigma_{i_u}(j_u)}$  which is substituted at that point. The leaves are  $\pm(\bar{q}_{-1} - \bar{q}_{-2})$ , which equals  $\pm \mathbf{d}_{-1}$  by Equation (5.9).

**Claim 5.3.5.** *The number of nodes in the recursion tree  $T_{\text{rec}}^{(i,j)}$  is at most  $2^i$ .*

*Proof.* Since  $T_{\text{rec}}^{(i,j)}$  is a binary tree, it suffices to prove that the height of  $T_{\text{rec}}^{(i,j)}$  is at most  $i$ . For each non-leaf node  $u$  of  $T_{\text{rec}}^{(i,j)}$ , let  $u_1$  be a child of  $u$ . According to Equation (5.13) and our labeling rules,  $i_{u_1}$  is a vertex in tetrahedron  $t_{i_u}$  and  $i_{u_1} < i_u$ . Let  $t_k$  be the tetrahedron such that  $t_k$  and  $t_{i_u}$  share a triangle face containing vertices  $\sigma_{i_u}(1), \sigma_{i_u}(2), \sigma_{i_u}(3)$ . Then in  $T_{\text{BFS}}$ ,  $t_k$  is the parent of  $t_{i_u}$ . Since  $i_{u_1} < i_u$ , in  $T_{\text{BFS}}$ , the depth of  $t_{i_{u_1}}$  is smaller than the depth of  $t_{i_u}$ . See Figure 5.3. It implies that the height of  $T_{\text{rec}}^{(i,j)}$  is at most the height of  $T_{\text{BFS}}$ , which is at most  $i$ .  $\square$

At each node  $u$  of the recursion tree  $T_{\text{rec}}^{(i,j)}$ , by applying Equation (5.13) with the left hand side being  $\bar{q}_{i_u} - \bar{q}_{\sigma_{i_u}(j_u)}$ , we introduce a term related to  $\mathbf{d}_{i_u}$ . For each non-root node  $u$ , denote  $\text{pa}(u)$  the parent of  $u$  in  $T_{\text{rec}}^{(i,j)}$ . Let  $P_u$  be the path from the root node to node  $u$  in

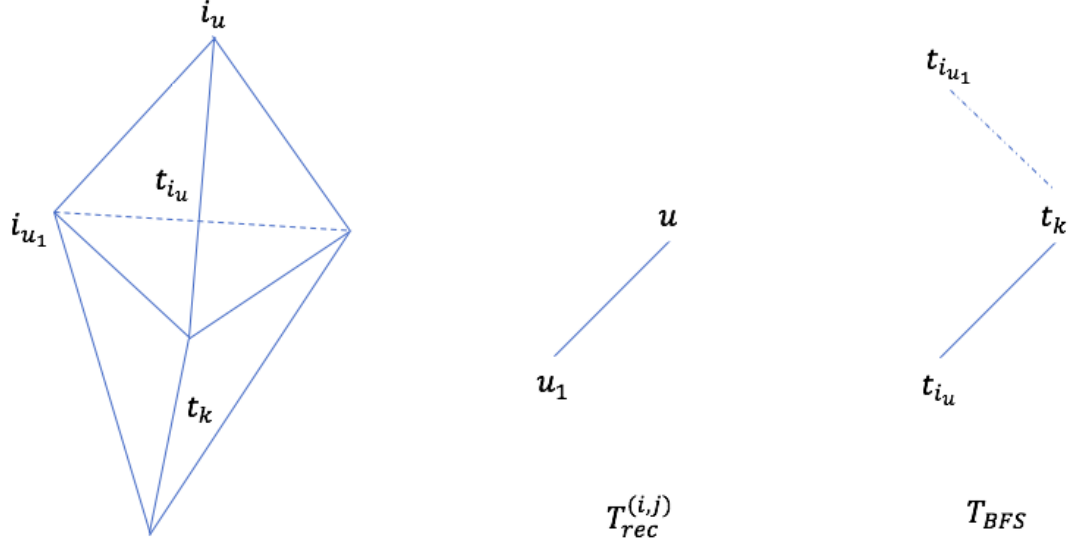


Figure 5.3: If  $u_1$  is a child of  $u$  in  $T_{rec}^{(i,j)}$ , then the depth of  $t_{i_{u_1}}$  is smaller than the depth of  $t_{i_u}$  in  $T_{BFS}$ .

$T_{rec}^{(i,j)}$ . For each non-root node  $u$ , define

$$\mathbf{M}_u \stackrel{\text{def}}{=} \prod_{v \in P_{pa(u)}} \frac{\mathbf{H}_{i_v, j_v}}{\det(\mathbf{F}_{i_v})}. \quad (5.14)$$

For root node  $r$ , define  $\mathbf{M}_r \stackrel{\text{def}}{=} \mathbf{I}$ . By recursively applying Equation (5.13), we get

$$\bar{\mathbf{q}}_i - \bar{\mathbf{q}}_{\sigma_i(j)} = \sum_{u \in T_{rec}^{(i,j)}} \text{sgn}(u) \mathbf{M}_u \mathbf{d}_{i_u},$$

where  $\text{sgn}(u) \in \{+1, -1\}$ . By the triangle inequality and the multiplicative inequality of 2-norm,

$$\|\bar{\mathbf{q}}_i - \bar{\mathbf{q}}_{\sigma_i(j)}\|_2 \leq \sum_{u \in T_{rec}^{(i,j)}} \|\mathbf{M}_u\|_2 \|\mathbf{d}_{i_u}\|_2. \quad (5.15)$$

We bound  $\|\mathbf{M}_u\|_2$  for each non-root node  $u$  by the following claim.

**Claim 5.3.6.** Let  $l_{\max}$  be the maximum edge length, and let  $V_{\min}$  be the minimum tetrahedron volume. Define  $c \stackrel{\text{def}}{=} \frac{l_{\max}^3}{6V_{\min}}$ , which is a constant by our assumption that all tetrahedrons have constant aspect ratios and all edge lengths are constant. For each non-root  $u$  in  $T_{\text{rec}}^{(i,j)}$ ,  $\|M_u\|_2 \leq c^i$ .

*Proof.* Let  $i = u_1 > u_2 > \dots > u_t > u_{t+1} = u$  be the vertices along the path  $P_u$ . By the proof of Claim 5.3.5, the length of  $P_u$  is at most  $i$ . Let  $i_k = i_{u_k}$  and  $j_k = j_{u_k}$  for each  $1 \leq k \leq t$ . By Equations (5.14), (5.11), and (5.12),

$$\begin{aligned} M_u &= \prod_{i_1 \geq i_l \geq i_t} \frac{\mathbf{h}_{i_l, j_l}(\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)})^\top}{(\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)})^\top \mathbf{h}_{i_l, j_l}} \\ &= \mathbf{h}_{i_1, j_1} \left( \prod_{i_1 \geq i_l \geq i_{t-1}} \frac{(\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)})^\top \mathbf{h}_{(i_{l+1}, j_{l+1})}}{(\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)})^\top \mathbf{h}_{(i_l, j_l)}} \right) \frac{(\mathbf{p}_{i_t} - \mathbf{p}_{\sigma_{i_t}(j_t)})^\top}{(\mathbf{p}_{i_t} - \mathbf{p}_{\sigma_{i_t}(j_t)})^\top \mathbf{h}_{(i_t, j_t)}}. \end{aligned}$$

Recall that from Claim 5.3.1,

$$(\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)})^\top \mathbf{h}_{(i_l, j_l)} = \det(\mathbf{F}_{i_l}),$$

which is equal to the volume of the tetrahedron generated by vectors  $\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(1)}$ ,  $\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(2)}$ , and  $\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(3)}$ . Similarly,  $(\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)})^\top \mathbf{h}_{(i_{l+1}, j_{l+1})}$  equals to the volume of the tetrahedron generated by vectors  $\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)}$ ,  $\mathbf{p}_{i_{l+1}} - \mathbf{p}_{\sigma_{i_{l+1}}((j_{l+1}+1)_3)}$ , and  $\mathbf{p}_{i_{l+1}} - \mathbf{p}_{\sigma_{i_{l+1}}((j_{l+1}+2)_3)}$ . Thus, by our definition  $c = l_{\max}^3/(6V_{\min}) \geq V_{\max}/V_{\min}$ ,

$$\left| \frac{(\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)})^\top \mathbf{h}_{(i_{l+1}, j_{l+1})}}{(\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)})^\top \mathbf{h}_{(i_l, j_l)}^{(j_l)}} \right| \leq c.$$

Thus,

$$\|M_u\|_2 \leq \left( \prod_{i_1 \geq i_l \geq i_{t-1}} \left| \frac{(\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)})^\top \mathbf{h}_{(i_{l+1}, j_{l+1})}}{(\mathbf{p}_{i_l} - \mathbf{p}_{\sigma_{i_l}(j_l)})^\top \mathbf{h}_{(i_l, j_l)}} \right| \right) \left\| \frac{\mathbf{h}_{(i_1, j_1)}(\mathbf{p}_{i_t} - \mathbf{p}_{\sigma_{i_t}(j_t)})^\top}{(\mathbf{p}_{i_t} - \mathbf{p}_{\sigma_{i_t}(j_t)})^\top \mathbf{h}_{(i_t, j_t)}} \right\|_2 = O(c^i).$$

□

Applying Claim 5.3.6 to Equation (5.15),

$$\|\bar{\mathbf{q}}_i - \bar{\mathbf{q}}_{\sigma_i(j)}\|_2 = O\left(c^i \sum_{u \in T_{\text{rec}}^{(i,j)}} \|\mathbf{d}_{i_u}\|_2\right).$$

By Claim 5.3.5 and the fact that  $-1 \leq i_u \leq i$  for each  $u \in T_{\text{rec}}^{(i,j)}$ ,

$$\|\bar{\mathbf{q}}_i - \bar{\mathbf{q}}_{\sigma_i(j)}\|_2 = O\left((2c)^i \sum_{-1 \leq i' \leq i} \|\mathbf{d}_{i'}\|_2\right). \quad (5.16)$$

### 5.3.4 Bounding the Quadratic Form $(\bar{\mathbf{q}})^\top \mathbf{M} \bar{\mathbf{q}}$

Note

$$(\bar{\mathbf{q}})^\top \mathbf{M} \bar{\mathbf{q}} \geq \sum_{\substack{-1 \leq i \leq n-3 \\ 1 \leq j \leq \min\{i+2, 3\}}} ((\bar{\mathbf{q}}_i - \bar{\mathbf{q}}_{\sigma_i(j)})^\top (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)}))^2.$$

By Equation (5.10),

$$(\bar{\mathbf{q}})^\top \mathbf{M} \bar{\mathbf{q}} \geq \sum_{\substack{-1 \leq i \leq n-3 \\ 1 \leq j \leq \min\{i+2, 3\}}} (\mathbf{d}_i^\top (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)}))^2.$$

By the fact that  $\mathbf{d}_{-1} = \bar{\mathbf{q}}_{-1} - \bar{\mathbf{q}}_{-2}$  is parallel to  $\mathbf{p}_{-1} - \mathbf{p}_{-2}$ ,

$$(\mathbf{d}_{-1}^\top (\mathbf{p}_{-1} - \mathbf{p}_{-2}))^2 = \Omega(\|\mathbf{d}_{-1}\|_2^2).$$

**Claim 5.3.7** (Lemma 3.7 of [44]). *Under the assumption: the angle between  $\mathbf{p}_0 - \mathbf{p}_{-1}$  and  $\mathbf{p}_0 - \mathbf{p}_{-2}$  is in the range  $[\theta, \pi - \theta]$  for some constant  $\theta$ . For any fixed unit vector  $\mathbf{y} \in \mathbb{R}^3$ ,*

$$\sum_{j \in \{-2, -1\}} (\mathbf{y}^\top (\mathbf{p}_0 - \mathbf{p}_j))^2 = \Theta(1).$$

**Claim 5.3.8.** *Under the assumption of constant edge lengths: for each  $1 \leq i \leq n - 3$ , the*

determinant of the matrix:

$$\det \left( \begin{pmatrix} \mathbf{p}_i - \mathbf{p}_{\sigma_i(1)} & \mathbf{p}_i - \mathbf{p}_{\sigma_i(2)} & \mathbf{p}_i - \mathbf{p}_{\sigma_i(3)} \end{pmatrix} \right) = \Theta(1).$$

Furthermore, for any fixed unit vector  $\mathbf{y} \in \mathbb{R}^3$ , for each  $1 \leq i \leq n-3$ ,

$$\sum_{1 \leq j \leq 3} (\mathbf{y}^\top (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)}))^2 = \Theta(1).$$

*Proof.* The determinant of matrix  $\begin{pmatrix} \mathbf{p}_i - \mathbf{p}_{\sigma_i(1)} & \mathbf{p}_i - \mathbf{p}_{\sigma_i(2)} & \mathbf{p}_i - \mathbf{p}_{\sigma_i(3)} \end{pmatrix}$  is the signed volume of tetrahedron  $t_i = \{i, \sigma_i(1), \sigma_i(2), \sigma_i(3)\}$ , which is constant by our assumption.

We claim that for any unit vector  $\mathbf{y}$ , there exists some  $j \in \{1, 2, 3\}$  such that  $\mathbf{y}^\top (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)}) = \Theta(1)$ . Assume by contradiction, for every  $j \in \{1, 2, 3\}$  we have  $\mathbf{y}^\top (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)}) = o(1)$ . This means that all three vectors  $\mathbf{p}_i - \mathbf{p}_{\sigma_i(1)}, \mathbf{p}_i - \mathbf{p}_{\sigma_i(2)}, \mathbf{p}_i - \mathbf{p}_{\sigma_i(3)}$  are between two 2D planes with distance  $o(1)$  and orthogonal to  $\mathbf{y}$ . This contradicts the assumption that the volume of tetrahedron  $t_i = \{i, \sigma_i(1), \sigma_i(2), \sigma_i(3)\}$  is constant.

Thus, we have

$$\sum_{1 \leq j \leq 3} (\mathbf{y}^\top (\mathbf{p}_i - \mathbf{p}_{\sigma_i(j)}))^2 = \Theta(1).$$

□

Therefore,

$$(\bar{\mathbf{q}})^\top \mathbf{M} \bar{\mathbf{q}} = \Omega \left( \sum_{-1 \leq i \leq n} \|\mathbf{d}_i\|_2^2 \right).$$

Applying the Cauchy-Schwarz inequality on the above equation and Equation (5.16) gives the path lemma.

## 5.4 Computing a $(B, r)$ -Hollowing of a Convex Edge-simple Truss

In this section, we present Algorithm 13 HOLLOW, which is used as a subroutine for line 6 of Algorithm 12 TRUSSSOLVER. The input of Algorithm 13 consists of a convex edge-simple 3D truss  $\mathcal{T} = \langle \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$ , a bounding box  $B$  and an integer parameter  $r$ . The algorithm outputs a  $(B, r)$ -hollowing  $\mathcal{H}$  of  $\mathcal{T}$ . We can check that the algorithm terminates in time  $O(|V|)$ . These together prove Lemma 5.1.3.

---

**Algorithm 13** HOLLOW( $\mathcal{T} = \langle \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle, B, r$ )

---

**Input:** a convex edge-simple 3D truss  $\mathcal{T} = \langle \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle$ , a bounding box  $B$ , a positive integer  $r \leq n/\alpha(\mathcal{T})^2$

**Output:** a  $(B, r)$ -hollowing of  $\mathcal{T}$

- 1: Compute  $r$ -division planes, which are planes orthogonal to each of the three directions of  $B$  such that these planes divide  $B$  into  $O(n/r)$  small cubes of side length  $O(r^{1/3})$  each.
  - 2:  $\mathcal{H} \leftarrow$  tetrahedrons intersecting an  $r$ -division plane.
  - 3:  $\mathcal{H} \leftarrow \mathcal{H} \cup$  tetrahedrons on the boundary of  $\mathcal{T}$ .
  - 4: Make  $\mathcal{H}$  stiffly-connected by adding a minimal number of tetrahedrons in  $\mathcal{T}$ .
  - 5: **return**  $\mathcal{H}$ .
- 

Given that each individual tetrahedron has constant volume and constant aspect ratio, the following observation converts counting the number of tetrahedrons in a truss intersecting a 2D plane into the intersection area of this truss and the plane.

**Observation 5.4.1.** *Let  $\mathcal{T}$  be a convex edge-simple 3D truss, and let  $P$  be a 2D plane. Let  $A(\mathcal{T} \cap P)$  be the intersection area of  $\mathcal{T}$  and  $P$ . Then the number of tetrahedrons in  $\mathcal{T}$  intersecting  $P$  is upper bounded by  $O(\max\{A(\mathcal{T} \cap P), 1\})$ .*

*Proof.* Since every individual tetrahedron in  $\mathcal{T}$  has constant volume and constant aspect ratio, a tetrahedron of  $\mathcal{T}$  intersects  $P$  only if all points of this tetrahedron is within some constant distance of  $P$ . The number of tetrahedrons of  $\mathcal{T}$  intersecting  $P$  can be upper bounded by the volume within some constant distance to  $\mathcal{T} \cap P$ . Thus, the number of tetrahedrons of  $\mathcal{T}$  intersecting  $P$  is at most  $O(\max\{A(\mathcal{T} \cap P), 1\})$ .  $\square$

### 5.4.1 Bounding the Size of $\mathcal{H}$

In this section, we show that  $\mathcal{H} = \text{HOLLOW}(\mathcal{T}, B, r)$ , computed by Algorithm 13, has a small size. That is,  $\mathcal{H}$  satisfies the first condition of the  $(B, r)$ -hollowing definition in Definition 5.1.2.

Note in Algorithm 13 line 1, the bounding box  $B$  is divided into  $O(n/r)$  small cubes of volume  $O(r)$  each. We call a small cube as a *region*. The tetrahedrons in a region are the tetrahedrons of  $\mathcal{H}$  which intersect a single small cube. A tetrahedron can appear in at most eight regions.

The following lemma upper bounds the number of tetrahedrons of  $\mathcal{H}$  in each region.

**Lemma 5.4.2.** *Given a convex edge-simple 3D truss  $\mathcal{T}$  of  $n$  vertices, a bounding box  $B$  and a positive integer  $r \leq n/\alpha(\mathcal{T})^2$ , let  $\mathcal{H} = \text{HOLLOW}(\mathcal{T}, B, r)$  returned by Algorithm 13. Then,  $\mathcal{H}$  has at most  $O(nr^{-1/3})$  tetrahedrons.*

Note the shortest side length of the bounding box of  $\mathcal{T}$  is at least  $n^{1/3}\alpha^{-2/3}$ . The requirement  $r \leq n/\alpha(\mathcal{T})^2$  guarantees that the shortest side length is at least  $r^{1/3}$  so that an  $r$ -division exists.

*Proof.* Note  $\mathcal{H}$  has  $O(n/r)$  regions. It suffices to show that each region of  $\mathcal{H}$  has at most  $O(r^{2/3})$  tetrahedrons.

Let  $R$  be a region of  $\mathcal{H}$ . A tetrahedron of  $\mathcal{H}$  belongs to region  $R$  if either this tetrahedron is within constant distance to the boundary of  $R$ , or this tetrahedron is within constant distance to the part of the boundary of  $\mathcal{T}$  that's contained in  $R$ . Since every tetrahedron of  $\mathcal{H}$  has constant volume and aspect ratio, the number of tetrahedrons within constant distance to the boundary of  $R$  is  $O(r^{2/3})$ . It remains to bound the number of tetrahedrons within constant distance to the boundary of  $\mathcal{T}$  that's contained in  $R$ .

Define  $S \stackrel{\text{def}}{=} \mathcal{T} \cap R$ . Since both  $\mathcal{T}$  and  $R$  are convex,  $S$  is convex. Let  $\text{Surf}(\cdot)$  the surface area of a shape. By Observation 5.4.1, the number of boundary tetrahedrons of  $\mathcal{T}$  contained in  $R$  is  $O(\text{Surf}(S))$ .

Let  $B_1$  be the smallest ball containing  $S$ . Since  $R$  is a cube, we have

$$\text{Surf}(B_1) = O(\text{Surf}(R)).$$

So it suffices to show  $\text{Surf}(S) \leq \text{Surf}(B_1)$ . We do so by giving a one-to-one mapping of every point from the surface of  $S$  onto  $B_1$ .

Let  $\phi : S \rightarrow B_1$  which maps each face of  $S$  to a subset of the surface of  $B_1$ , defined as follows. Consider a face of  $S$ , say  $f$ , with vertices  $v_1, \dots, v_k$  in a clockwise order. Let  $P_f$  be the plane containing  $f$ .  $P_f$  cuts  $B_1$  into two parts, let  $B'_1$  be the part of the smaller volume (break a tie arbitrarily), aka the sphere cap generated by the plane  $p_f$ .

For each  $1 \leq i \leq k$ , let  $P_i$  be the plane orthogonal to  $P_f$  that passes through  $v_i$  and  $v_{i+1}$  (if  $i = k$ , then the intersection line is  $(v_k, v_1)$ ). Let  $u_i$  be the point of intersection of the surface of  $B_1$  with the planes  $P_{i-1}$  and  $P_i$  (if  $i = 1$ , then  $u_1$  is the intersect vertex of  $P_1, P_k$  and the surface of  $B'_1$ ).

We define  $\phi(f)$  to be the surface of  $B'_1$  enclosed by  $(u_1, \dots, u_k, u_1)$ . See Figure 5.4 for an example.

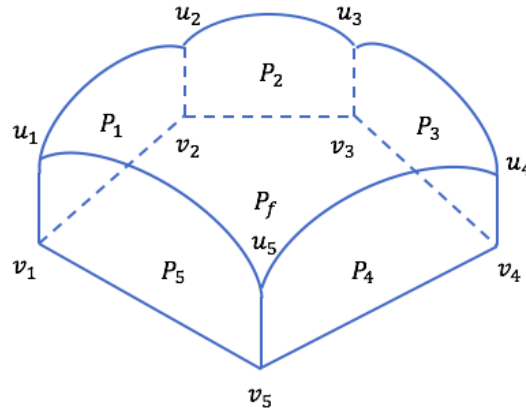


Figure 5.4: An example of  $\phi(f)$ . Face  $f$  has five vertices  $v_1, \dots, v_5$ .  $\phi(f)$  is the surface of  $B'_1$  enclosed by  $(u_1, \dots, u_5, u_1)$ .



For each face  $f$  of  $S$ , the orthogonal projection of  $\phi(f)$  onto the plane  $P_f$  is  $f$ . Thus

$$\text{AREA}(\phi(f)) \geq \text{AREA}(f).$$

In addition, since  $S$  is convex, for any two distinct faces  $f_1$  and  $f_2$ ,  $\phi(f_1)$  and  $\phi(f_2)$  are also disjoint. So we have

$$\text{Surf}(S) = \sum_{f \in S} \text{AREA}(f) \leq \sum_{f \in S} \text{AREA}(\phi(f)) \leq \text{Surf}(B_1).$$

Combining this with  $\text{Surf}(B_1) \leq O(\text{Surf}(S)) \leq O(r^{2/3})$  then completes the proof.  $\square$

#### 5.4.2 Bounding the Number of Tetrahedrons in $\mathcal{H}$ Intersecting with a Plane

In this section, we show that  $\mathcal{H} = \text{HOLLOW}(\mathcal{T}, B, r)$ , computed by Algorithm 13, has a small overlap with any plane whose normal vector has an angle between  $(0, \pi/2)$  with the longest direction of  $B$ . That is,  $\mathcal{H}$  satisfies the second condition of the  $(B, r)$ -hollowing definition in Definition 5.1.2.

**Lemma 5.4.3.** *Given a convex edge-simple 3D truss  $\mathcal{T}$  of  $n$  vertices, a bounding box  $B$  and a positive integer  $r \leq n/\alpha(\mathcal{T})^2$ , let  $\mathcal{H} = \text{HOLLOW}(\mathcal{T}, B, r)$  returned by Algorithm 13. Let  $\mathbf{d} \in \mathbb{R}^3$  be a unit vector such that the angle between  $\mathbf{d}$  and let the angles with the three directions of the box (normals to its faces) be  $\theta_x, \theta_y, \theta_z \in (\theta, \pi/2)$ , for some  $\theta > 0$ . Then, the number of tetrahedrons in  $\mathcal{H}$  which intersect any plane  $P$  orthogonal to  $\mathbf{d}$  is at most*

$$O\left(\frac{n^{2/3}}{\alpha(\mathcal{T})^{1/3} r^{1/3} \cos^2 \theta}\right).$$

Without loss of generality, we assume that the bounding box  $B$  is *axis-parallel*, that is, the sides of  $B$  are parallel to the three axes: the  $x$  axis, the  $y$  axis and the  $z$ -axis. We say an axis-parallel box has side lengths  $a, b, c > 0$ , if the sides parallel to the  $x$ -axis have

length  $a$ , the sides parallel to the  $y$ -axis have length  $b$ , and the sides parallel to the  $z$ -axis have length  $c$ , respectively.

To prove Lemma 5.4.3, we need the following claim, which bounds the intersection area of a 2D plane and a 3D box.

**Claim 5.4.4.** *Let  $B$  be a 3D axis-parallel box of side lengths  $a, b, c > 0$ . Let  $\mathbf{d} \in \mathbb{R}^3$  be a unit vector such that the angle between  $\mathbf{d}$  and the  $x$ -axis (the  $y$ -axis, and the  $z$ -axis) is  $\theta_x$  ( $\theta_y, \theta_z$ , respectively). Suppose  $\theta_x, \theta_y, \theta_z \in (0, \pi/2)$ . Then the intersection area  $S$  of any 2D plane  $P$  orthogonal to  $\mathbf{d}$  and the 3D box  $B$  satisfies*

$$S \leq \min \left\{ \frac{bc}{\cos \theta_x}, \frac{ac}{\cos \theta_y}, \frac{ab}{\cos \theta_z} \right\}.$$

*Proof.* Since the three terms in the right hand side are symmetric, we only prove  $S$  can be upper bounded by the first term and the other two follow in a similar way.

Let  $B_{x1}, B_{x2}$  be the two faces of  $B$  which are orthogonal to the  $x$ -axis, without loss of generality, assume  $B_{x1}$  has a smaller  $x$ -coordinate.

If  $P$  intersects neither  $B_{x1}$  nor  $B_{x2}$ , then the volume of  $B$  is equal to  $Sa \cos \theta_x$ .

If  $P$  intersects  $B_{x1}$  say with line  $(K, L)$ , see Figure 5.5, then we draw a line going through point  $K$  and parallel to the  $x$ -axis, which intersects face  $B_{x2}$  at point  $M$ , similarly we draw a line going through point  $L$  and parallel to the  $x$ -axis, which intersects face  $B_{x2}$  at point  $N$ . We cut the box  $B$  by the plane  $KMLN$ , see Figure 5.5. Note that the volume of the convex hull of  $(K, G, E, F, G, H, I, J)$ , the right one in Figure 5.5, equals to  $Sa \cos \theta_x$ , which is smaller than the volume of  $B$ .

Similarly, if  $P$  intersects  $B_{x2}$ , then we can draw two lines parallel to the  $x$ -axis and going through the two intersection points respectively and get a shape of volume  $Sa \cos \theta_x$  smaller than the volume of  $B$ . Note that the intersection between  $P$  and  $B_{x2}$  cannot coincide with  $(M, N)$ , otherwise the angle  $\theta_x = \pi/2$ . We can check that the shape we get

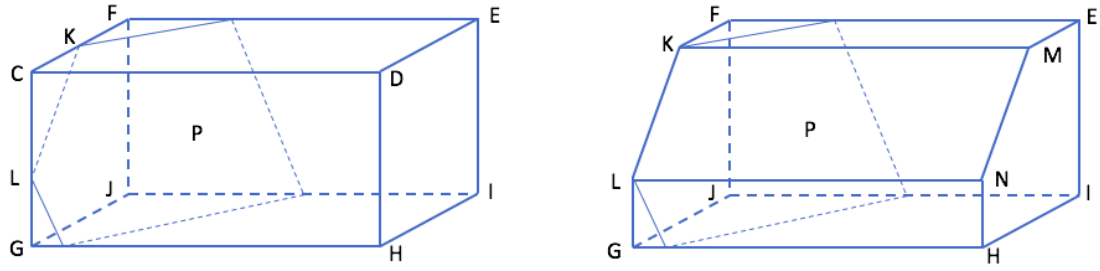


Figure 5.5: Plane  $P$  intersects face  $B_{x1}$  with line  $(K, L)$ . We draw a line going through point  $K$  and parallel to the  $x$ -axis, which intersects face  $B_{x2}$  at point  $M$ , similarly we draw a line going through point  $L$  and parallel to the  $x$ -axis, which intersects face  $B_{x2}$  at point  $N$ .

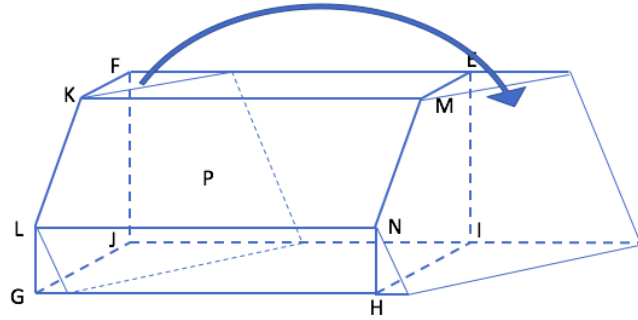


Figure 5.6: Cut the shape along the plane  $P$ , and shift the left part along the  $x$ -axis and glue the two faces  $KFLJGL$  and  $MEIHN$ .

must contain  $P \cap B$ , and the two faces of the shape which are orthogonal to the  $x$ -axis are congruent. If we cut the shape along the plane  $P$ , then we can shift the left part along the  $x$ -axis and glue the two faces which are orthogonal to the  $x$ -axis, and get a parallelepiped which has  $P \cap B$  as a face. Figure 5.6 shows the parallelepiped we get from the example of Figure 5.5. Thus,

$$S \cdot a \cdot \cos \theta_x \leq abc.$$

That is,  $S \leq bc / \cos \theta_x$ . This completes the proof.  $\square$

Now we prove Lemma 5.4.3.

*Proof of Lemma 5.4.3.* Let  $t$  be the number of regions (that is, small cubes of the hollowing) of  $\mathcal{H}$  which intersect a 2D plane  $P$ , and let  $m$  be the maximum number of tetrahedrons of a single region of  $\mathcal{H}$  which intersect the plane  $P$ . The number of tetrahedrons in the hollowing  $\mathcal{H}$  which intersect  $P$  can be upper bounded by  $tm$ .

We first bound  $t$ , the number of regions of  $\mathcal{H}$  which intersect the plane  $P$ . A cube region intersects the plane  $P$  only if all its points are within a distance  $\sqrt{3}r^{1/3}$  of  $P$ . We put two planes, say  $P_1, P_2$ , which are parallel to the plane  $P$ , above and below  $P$  with distance  $\sqrt{3}r^{1/3}$  to  $P$ . All cube regions which intersect the plane  $P$  must be within the two planes  $P_1$  and  $P_2$ . By Claim 5.4.4, the volume of the box  $B$  between the two planes  $P_1, P_2$  is at most

$$\frac{bc}{\cos \theta} \cdot 2\sqrt{3}r^{1/3}.$$

Since each cube has volume at most  $r$ , we can bound the number of cube regions which intersect the plane  $P$

$$t \leq \frac{2\sqrt{3}bc}{r^{2/3} \cos \theta}.$$

We then bound  $m$ , the maximum number of tetrahedrons of a single region of  $\mathcal{H}$  which intersect the plane  $P$ . Consider the tetrahedrons of  $\mathcal{H}$  in this single cube region which intersect a some hollowing plane. These tetrahedrons are within constant distance of an  $r^{1/3} \times r^{1/3}$  square, which is on a plane orthogonal to one of the direction of the bonding box. By Observation 5.4.1 and Claim 5.4.4, the number of these tetrahedrons which intersect the plane  $P$  is at most  $O(r^{1/3}/\cos \theta)$ . Thus, we can bound the maximum number of tetrahedrons of a single cube region which intersect the plane  $P$

$$m \leq O\left(\frac{r^{1/3}}{\cos \theta}\right).$$

Therefore, the number of tetrahedrons of the hollowing  $\mathcal{H}$  which intersect the plane  $P$  is at

most

$$tm \leq O\left(\frac{bc}{r^{1/3} \cos^2 \theta}\right).$$

Note that we have  $a = c \cdot \alpha(\mathcal{T})$  and  $a \geq b \geq c$ . The above number is upper bounded by

$$O\left(\frac{n^{2/3}}{\alpha(\mathcal{T})^{1/3} r^{1/3} \cos^2 \theta}\right).$$

□

### 5.4.3 Bounding the Relative Condition Number of $\mathcal{H}$

In this section, we show that the condition number of  $\mathbf{A}_{\mathcal{H}}$  and  $\text{SC}[\mathbf{A}_{\mathcal{T}}]_U$  is small. That is,  $\mathcal{H}$  satisfies the third condition of the  $(B, r)$ -hollowing definition in Definition 5.1.2.

Since  $\mathcal{T}$  is convex, each region of  $\mathcal{H}$  is connected. Lemma 5.1.1 implies that  $\mathbf{A}_{\mathcal{H}}$  and  $\text{SC}[\mathbf{A}_{\mathcal{T}}]_U$  has the same null space. By Lemma 5.1.1, the smallest nonzero eigenvalue of  $\mathbf{A}_{\mathcal{H}}$  is lower bounded by the diameter and the size of each region of  $\mathcal{H}$ .

**Lemma 5.4.5.** *Let  $\mathbf{A}_{\mathcal{T}}$  be the truss stiffness matrix of a convex edge-simple 3D truss  $\mathcal{T}$ . Given a positive integer  $r$  and a bounding box  $B$ , let  $\mathcal{H} = \text{HOLLOW}(\mathcal{T}, B, r)$  be returned by Algorithm 13. Let  $\mathbf{A}_{\mathcal{H}}$  be the associated truss stiffness matrix of  $\mathcal{H}$ . Then,*

$$\mathbf{A}_{\mathcal{H}} \preceq \text{SC}[\mathbf{A}_{\mathcal{T}}]_U \preceq O(r^2) \mathbf{A}_{\mathcal{H}},$$

where  $U$  consists of all vertices in  $\mathcal{H}$ .

*Proof.* The first inequality is equivalent to:  $\text{SC}[\mathbf{A}_{\mathcal{T}}]_U - \mathbf{A}_{\mathcal{H}}$  is a symmetric PSD matrix. Let  $\mathcal{T}'$  be the truss obtained by removing all the edges of  $\mathcal{T}$  whose two endpoints are both in  $U$ . Let  $\mathbf{A}_{\mathcal{T}'}$  be the associated truss stiffness matrix of  $\mathcal{T}'$ . We can check that

$$\text{SC}[\mathbf{A}_{\mathcal{T}}]_U - \mathbf{A}_{\mathcal{H}} = \text{SC}[\mathbf{A}_{\mathcal{T}'}]_U.$$

By Fact 2.1.3,  $\text{SC}[\mathbf{A}_{\mathcal{T}}]_U - \mathbf{A}_{\mathcal{H}}$  is a symmetric PSD matrix.

It remains to prove the second inequality. Note the planes of Algorithm 13 line 1 divide  $\mathcal{T}$  into small regions. Let  $\mathcal{T}_i$  denote the subgraph induced by  $\mathcal{T}$  on the  $i$ th region, and let  $\mathcal{T}_{C_i}$  denote the subgraph induced by  $\mathcal{T}$  on the boundary of the  $i$ th region. Let  $\mathbf{A}_{\mathcal{T}_i}, \mathbf{A}_{\mathcal{T}_{C_i}}$  be the associated truss stiffness matrices of  $\mathcal{T}_i$  and  $\mathcal{T}_{C_i}$  respectively. Let  $U_i$  denote the boundary vertices of the  $i$ th region. Let  $\text{SC}[\mathbf{A}_{\mathcal{T}_i}]_{U_i}$  be the Schur complement of  $\mathbf{A}_{\mathcal{T}_i}$  w.r.t. to  $U_i$ .

Since  $\mathcal{T}$  is a convex edge-simple 3D truss, each vertex in  $\mathcal{T}_i$  has constant degree and each edge has constant length and elasticity parameter. It implies  $\lambda_{\max}(\mathbf{A}_{\mathcal{T}_i}) = O(1)$ . By Fact 2.1.3,  $\lambda_{\max}(\text{SC}[\mathbf{A}_{\mathcal{T}_i}]_{U_i}) = O(1)$ .

According to Algorithm 13 from line 2 to line 4, in each  $\mathcal{T}_{C_i}$ , the tetrahedrons are arranged in simplicial complex and  $\mathcal{T}_{C_i}$  is connected. By Lemma 5.1.1, the null spaces of  $\mathbf{A}_{\mathcal{T}_{C_i}}$  and the null space of  $\text{SC}[\mathbf{A}_{\mathcal{T}_i}]_{U_i}$  are the same. Besides, each  $\mathcal{T}_{C_i}$  has  $O(r^{2/3})$  vertices and diameter  $O(r^{1/3})$ , by Lemma 5.4.2. Applying Lemma 5.1.1 gives:

$$\lambda_{\min}(\mathbf{A}_{\mathcal{T}_{C_i}}) = \Omega\left(\frac{1}{r^2}\right).$$

This implies:

$$\text{SC}[\mathbf{A}_{\mathcal{T}_i}]_{U_i} \preceq O(r^2) \mathbf{A}_{\mathcal{T}_{C_i}}.$$

Note each edge of  $\mathcal{T}$  only appears in a constant number of regions. Thus,

$$\text{SC}[\mathbf{A}_{\mathcal{T}}]_U \preceq \sum_i \text{SC}[\mathbf{A}_{\mathcal{T}_i}]_{U_i} \preceq O(r^2) \sum_i \mathbf{A}_{\mathcal{T}_{C_i}} \preceq O(r^2) \mathbf{A}_{\mathcal{H}}.$$

This completes the proof. □

## 5.5 Nested Dissection

In this section, we present our solver for the constructed preconditioners, Algorithm 15 CONVEXTRUSSUNIONND, which is used as a subroutine in line 8 of Algorithm 12 TRUSS-

SOLVER.

We first prove Lemma 2.5.2 in Section 5.5.1. Then, in Section 5.5.2, we use Lemma 2.5.2 to analyze the performance of Algorithm 15, which proves Lemma 5.1.5.

### 5.5.1 Proof of Lemma 2.5.2

We restate Lemma 2.5.2 with running time phrased in terms of the matrix multiplication exponent  $\omega$ .

**Lemma 2.5.2.** *Suppose we have a recursive separator decomposition of a simplicial complex with  $n$  bounded aspect ratio tetrahedrons such that:*

1. *the number of leaves, and hence total number of recursive calls, is at most  $n^\alpha$ .*
2. *each leaf (bottom layer partition) has at most  $n^\beta$  tetrahedrons.*
3. *each top separator has size at most  $n^\gamma$ .*

*Then we can find an exact Cholesky factorization of the associated stiffness matrix with multiplication count  $O(n^{\alpha+2\omega\beta/3} + n^{\alpha+\gamma\omega} \log n)$  and fill-in size  $O(n^{\alpha+\frac{4}{3}\beta} + n^{\alpha+2\gamma})$ .*

Nested dissection according to the separator decomposition stated in Lemma 2.5.2 has three parts of cost:

1. Inverting leaf components:
  - (a) the cost only associated with vertices not belonging to top-level separators;
  - (b) the cost associated with top-level separators.
2. Inverting top-level separators.

We first analyze the cost of inverting top-level separators. Note that top-level separators are disjoint. After eliminating all leaf components, we get a *layered graph*. That is, its vertices can be partitioned into  $V_1 \cup \dots \cup V_l$  for some positive integer  $l$  such that there is

no edge between  $V_i$  and  $V_j$  for  $|i - j| > 1$ . Algorithm 14 gives a numbering for a layered graph, and Claim 5.5.1 analyzes its performance.

---

**Algorithm 14** LAYEREDGRAPHND( $G = (V_1 \cup \dots \cup V_l, E)$ )

---

**Input:** a layered graph  $G$  with  $l$  layers

**Output:** an elimination ordering for vertices in  $G$

- 1: **if**  $l = 1$  **then**
  - 2:   number the vertices in  $V_1$  arbitrarily and return this numbering.
  - 3: **end if**
  - 4: Label  $V_{\lfloor l/2 \rfloor}$  with the highest possible numbers.
  - 5: **return** the numbering of LAYEREDGRAPHND( $G[V_1 \cup \dots \cup V_{\lfloor l/2 \rfloor - 1}]$ ) and LAYEREDGRAPHND( $G[V_{\lfloor l/2 \rfloor + 1} \cup \dots \cup V_l]$ ), and  $V_{\lfloor l/2 \rfloor}$ .
- 

**Claim 5.5.1.** *Let  $G$  be a layered graph with  $l$  layers of at most  $s$  vertices each. Algorithm 14 LAYEREDGRAPHND returns a numbering such that, Gaussian elimination according to this order has fill-in size  $O(s^2l)$  and multiplication count  $O(s^\omega l)$ .*

*Proof.* We follow the proofs of [48]. We first bound the fill-in size. Consider a recursion of Algorithm 14 on a graph with  $n$  vertices, let  $f(n)$  denote the maximum number of fill-in edges whose lower numbered endpoint is numbered by this recursion. Suppose this recursion deals with layers  $V_s, V_{s+1}, \dots, V_t$ . The algorithm numbers the vertices in the middle layer, that is,  $V_{\lfloor (s+t)/2 \rfloor}$ , which can be viewed as a separator whose removal separates the graph into two disjoint and balanced parts. The fill-in edges whose lower numbered endpoint is in  $V_{\lfloor (s+t)/2 \rfloor}$  consists of the following edges: (1) edges whose both endpoints are in  $V_{\lfloor (s+t)/2 \rfloor}$ ; and (2) edges whose one endpoint is in  $V_{\lfloor (s+t)/2 \rfloor}$  and the other is in  $V_{s-1}$  (if exists) or  $V_{t+1}$  (if exists). Since each layer has at most  $s$  vertices,

$$f(n) \leq \begin{cases} O(s^2), & \text{if } n = O(s) \\ f(n_1) + f(n_2) + 3s^2, & \text{otherwise} \end{cases}$$

Note  $n_2 \leq n_1 \leq n_2 + s$ . Thus, the total fill-in size is  $O(s^2l)$ .

We then bound the multiplication count. For a recursion of Algorithm 14 on a graph



with  $n$  vertices, let  $g(n)$  denote the maximum multiplication count associated with vertices in this graph which are going to be numbered in this recursion. Similar to the analysis of the fill-in, we have

$$g(n) \leq \begin{cases} O(s^\omega), & \text{if } n = O(s) \\ g(n_1) + g(n_2) + 3s^\omega, & \text{otherwise} \end{cases}$$

Thus, the total multiplication count is  $O(s^\omega l)$ . □

Using Algorithm 14 as a subroutine, we prove Lemma 2.5.2.

*Proof of Lemma 2.5.2.* We first bound the fill-in size. Note each leaf component has  $O(n^\beta)$  vertices in which  $O(n^\gamma)$  vertices are in some top-level separators and are labeled by higher numbers. By the result in [48] and [49], the fill-in introduced by inverting each leaf component is

$$O(n^{4\beta/3} + n^{2\beta/3+\gamma}).$$

There are totally  $O(n^\alpha)$  leaf components. Thus, the fill-in size introduced by inverting all leaf components is  $O(n^{\alpha+4\beta/3} + n^{\alpha+2\beta/3+\gamma})$ .

By Claim 5.5.1, the fill-in size of inverting top-level separators is  $O(n^{\alpha+2\gamma})$ . Thus, the total fill-in size is

$$O(\alpha + n^{4\beta/3} + n^{\alpha+2\beta/3+\gamma} + n^{\alpha+2\gamma}) = O(n^{\alpha+4\beta/3} + n^{\alpha+2\gamma}).$$

Here, we use Cauchy-Schwarz inequality to drop the second term.

We then bound the multiplication count. By [49], when inverting each leaf component, the multiplication count only associated with vertices not belonging to top-level separators is  $O(n^{2\omega\beta/3})$ . By Claim 5.5.1, the multiplication count of inverting top-level separators is  $O(n^{\alpha+\omega\gamma})$ .

It remains to upper bound the multiplication count associated with top-level separators

when inverting a leaf component. We adapt the analysis in [48].

We prove a more general result. Consider a tetrahedral mesh  $G$  with  $n$  vertices, in which  $s$  vertices are in a top-level separator and have been labeled with higher numbers. Let  $g(n, s)$  be the multiplication count associated these  $s$  top-level vertices when running nested dissection on  $G$ .

Nested dissection finds a separator of size  $n^{2/3}$  for  $G$ . The multiplication count associated with the  $s$  top-level vertices when inverting this separator can be upper bounded by  $O((n^{2/3} + s)^\omega)$ . This gives the following recursion:

$$\begin{aligned} g(n, s) &\leq O((n^{2/3} + s)^\omega) + \max_{n_i, s_i} \left\{ \sum_{1 \leq i \leq 2} g(n_i, s_i) \right\} \\ &\leq O(2^\omega n^{2\omega/3} + 2^\omega s^\omega) + \max_{n_i, s_i} \left\{ \sum_{1 \leq i \leq 2} g(n_i, s_i) \right\}. \end{aligned}$$

Here, the maximum is taken over

$$\begin{aligned} s_1 + s_2 &\leq s + n^{2/3} \\ n_1 + n_2 &\leq n + n^{2/3} \\ n_1, n_2 &\leq \frac{4}{5}n + n^{2/3}. \end{aligned}$$

We can compute that

$$g(n, s) = O(n^{2\omega/3} + s^\omega \log n).$$

Note each leaf component is incident to at most two top-level separators. Thus, the multiplication count associated with the top-level vertices when inverting this leaf component is

$$O(n^{2\omega\beta/3} + n^{\gamma\omega} \log n).$$

Combining the other two parts of cost, and the fact that there are totally  $n^\alpha$  leaf components,

the total multiplication count is

$$O\left(n^{\alpha+2\omega\beta/3} + n^{\alpha+\gamma\omega} \log n\right).$$

If we replace  $\omega$  by 3, then we can drop the  $\log n$  factor in the above equation, which gives the result in the statement.  $\square$

### 5.5.2 Proof of Lemma 5.1.5

We now present Algorithm 15 CONVEXTRUSSUNIONND. The input consists of a edge-simple 3D truss  $\mathcal{T}$  which is a union of  $k$  convex edge-simple trusses, a bounding box for each convex edge-simple truss, the index subset of small-aspect-ratio trusses and  $(B_i, r_i)$ -hollowings for each small-aspect-ratio truss. The output is an elimination ordering for the union of the hollowings of small-aspect-ratio trusses and the large-aspect-ratio trusses. This proves Lemma 5.1.5.

We first prove that there exists a good direction  $\mathbf{d}$  such that: the angle between  $\mathbf{d}$  and the longest direction of each bounding box is in a proper range.

**Lemma 5.5.2.** *Let  $k \geq 2$  and  $\mathbf{d}_1, \dots, \mathbf{d}_k \in \mathbb{R}^3$  be  $k$  unit vectors. Then there exists a unit vector  $\mathbf{d} \in \mathbb{R}^3$  such that*

$$\frac{1}{10k} \leq |\mathbf{d}^\top \mathbf{d}_i| \leq 1 - \frac{1}{10k}, \forall 1 \leq i \leq k. \quad (5.17)$$

*Proof.* We pick a unit vector  $\mathbf{d}$  uniformly at random. For any fixed  $i$ ,

$$\Pr\left(\frac{1}{10k} \leq |\mathbf{d}^\top \mathbf{d}_i| \leq 1 - \frac{1}{10k}\right) = 2 \cdot \frac{\text{vol}(1/10k) - \text{vol}(1 - 1/10k)}{V}.$$

Here,  $\text{vol}(x)$  is the volume of a cap of a 3D unit ball with height  $1 - x$ , and  $V$  is the volume of a 3D unit ball. We have  $\text{vol}(x) = \frac{1}{6}\pi(1 - x^2)(3(1 - x^4) + (1 - x^2)^2)$  and  $V = \frac{4}{3}\pi$ .

---

**Algorithm 15** CONVEXTRUSSUNIONND( $\mathcal{T} = \langle \{\mathbf{p}_i\}_{i \in V}, T, E, \gamma \rangle, \mathcal{B}, \mathcal{I}, \mathcal{H}, l$ )

---

**Input:** a edge-simple 3D truss  $\mathcal{T}$  which is the union of  $k$  convex edge-simple trusses

$\mathcal{T}_1, \dots, \mathcal{T}_k$ ,

$\mathcal{B} = \{B_1, \dots, B_k\}$  in which  $B_i$  is a bounding box of  $\mathcal{T}_i$  constructed from Line 2 of Algorithm 12,

Index set  $\mathcal{I}$  of large aspect ratio complexes constructed from Line 4 of Algorithm 12,

Hollowings of each  $\mathcal{T}_i$  with  $i \in \mathcal{I}$  constructed from Line 6 of Algorithm 12,

$l$ , the number of top-level partitions.

**Output:** an elimination ordering for vertices in  $(\cup_{i \in \mathcal{I}} \mathcal{H}_i) \cup (\cup_{i \notin \mathcal{I}} \mathcal{T}_i)$

- 1: For each  $1 \leq i \leq k$ ,  $\mathbf{d}_i \leftarrow$  the direction of the longest sides of  $B_i$ .
  - 2: Compute a unit vector  $\mathbf{d}$  such that  $\frac{1}{10k} \leq |\mathbf{d}^\top \mathbf{d}_i| \leq 1 - \frac{1}{10k}, \forall 1 \leq i \leq k$ .
  - 3: Compute separator planes  $P_1, \dots, P_l$ , which are planes orthogonal to  $\mathbf{d}$  and dividing  $\mathcal{T}$  into  $l + 1$  parts  $Q_1, \dots, Q_{l+1}$  of  $O(nl^{-1})$  tetrahedrons each.
  - 4: For each  $1 \leq j \leq l$ ,  $S_j \leftarrow$  tetrahedrons in  $(\cup_{i \in \mathcal{I}} \mathcal{H}_i) \cup (\cup_{i \notin \mathcal{I}} \mathcal{T}_i)$  which intersect plane  $P_j$ .
  - 5:  $Q_1 \leftarrow Q_1 \cup S_1, Q_{l+1} \leftarrow Q_{l+1} \cup S_l$ .
  - 6: For each  $2 \leq j \leq l$ ,  $Q_j \leftarrow Q_j \cup S_{j-1} \cup S_j$ .
  - 7: LAYEREDGRAPHND( $\mathcal{T}[S_1 \cup \dots \cup S_l]$ ) with highest numbers.
  - 8: Run nested dissection with MT-separators<sup>1</sup> for each  $Q_j$  to number its unnumbered vertices.
  - 9: **return** the elimination ordering of vertices in  $(\cup_{i \in \mathcal{I}} \mathcal{H}_i) \cup (\cup_{i \notin \mathcal{I}} \mathcal{T}_i)$ .
- 

Plugging these volumes into the above equation,

$$\Pr \left( \frac{1}{10k} \leq |\mathbf{d}^\top \mathbf{d}_i| \leq 1 - \frac{1}{10k} \right) = \frac{1}{2} \left( \left( 1 - \frac{1}{10k} \right)^2 \left( 2 + \frac{1}{10k} \right) - \left( \frac{1}{10k} \right)^2 \left( 3 - \frac{1}{10k} \right) \right).$$

Take the opposite:

$$\Pr \left( |\mathbf{d}^\top \mathbf{d}_i| < \frac{1}{10k} \text{ or } |\mathbf{d}^\top \mathbf{d}_i| > 1 - \frac{1}{10k} \right) \leq \frac{3}{20k} - \frac{1}{2000k^3} + \frac{3}{200k^2}.$$

By union bound,

$$\Pr \left( \exists i \text{ s.t. } |\mathbf{d}^\top \mathbf{d}_i| < \frac{1}{10k} \text{ or } |\mathbf{d}^\top \mathbf{d}_i| > 1 - \frac{1}{10k} \right) \leq k \left( \frac{3}{20k} - \frac{1}{2000k^3} + \frac{3}{200k^2} \right) \leq \frac{1}{5}.$$

Thus,

$$\Pr \left( \forall i, \frac{1}{10k} \leq |\mathbf{d}^\top \mathbf{d}_i| \leq 1 - \frac{1}{10k} \right) \geq \frac{4}{5}.$$

This implies there exists a  $\mathbf{d}$  as desired.  $\square$

We independently pick  $O(\log n)$  unit vectors uniformly at random. By a Chernoff bound, we can find a direction  $\mathbf{d}$  satisfying Equation (5.17) with high probability.

Recall that Lemma 5.4.3 states that: any 2D plane  $P$  orthogonal to  $\mathbf{d}$  intersects a small number of tetrahedrons in a  $(B, r)$ -hollowing of a convex edge-simple 3D truss. Lemma 5.1.4 bounds the number of tetrahedrons in a convex edge-simple 3D truss which intersect a single 2D plane  $P$  orthogonal to  $\mathbf{d}$ . We restate it in the following, which can be proved by combining Observation 5.4.1 and Claim 5.4.4.

**Lemma 5.1.4.** *Let  $\mathcal{T}$  be a convex edge-simple 3D truss of  $n$  vertices, and let  $B$  be a bounding box of  $\mathcal{T}$ . Let  $\mathbf{d} \in \mathbb{R}^3$  be a unit vector such that the angle between  $\mathbf{d}$  and the longest direction of  $B$  is  $\theta \neq \pi/2$ . Then any plane  $P$  orthogonal to  $\mathbf{d}$  intersects  $\mathcal{T}$  in at most  $O(n^{2/3} \alpha(\mathcal{T})^{-1/3} \cos^{-1} \theta)$  tetrahedrons.*

Lemma 5.4.3 and Lemma 5.1.4 together imply that: the top-level separator  $S_1 \cup \dots \cup S_l$  computed in Algorithm 15 has a small size. This, together with nested dissection in [49], lets us prove that Algorithm 15 outputs an elimination ordering with a small fill-in size and multiplication count.

**Lemma 5.5.3.** *Given a edge-simple 3D truss  $\mathcal{T}$  of  $n$  vertices which is a union of  $k$  convex edge-simple trusses with  $n_i$  vertices each, running Algorithm 12,  $\text{TRUSSSOLVER}(\mathcal{T}, \mathbf{f}, \epsilon, c_\alpha, c_r)$ , with Line 8 replaced by Algorithm 15,  $\text{CONVEXTRUSSUNIONND}(\mathcal{T}, \mathcal{B}, \mathcal{I}, \mathcal{H}, l)$ , leads to performance in terms of  $n$  that is optimized by setting*

$$c_\alpha = c_r \leq \frac{1}{3}$$

*in Line 4 of Algorithm 12,  $\text{TRUSSSOLVER}$ . In terms of  $c_r$ , the hollowing parameter, and  $l$ , the number of top-level separators, this gives an elimination ordering with fill-in size at most*

$$O(n^{4/3} l^{-1/3} + k^{14/3+2c_r/3} n^{4/3-2c_r/3} l),$$

that can be computed in time

$$O(n^{2\omega/3}l^{-2\omega/3+1} + k^{7\omega/3+\omega c_r/3}n^{2\omega/3-c_r\omega/3}l\log n),$$

where  $\omega$  is the matrix multiplication exponent.

*Proof.* We apply Lemma 2.5.2. According to Algorithm 15 line 2, for each  $i$ , the angle between the longest direction of the  $i$ th bounding box and  $\mathbf{d}$  has cosine value in  $[1/10k, 1 - 1/10k]$ .

We first upper bound the number of vertices in each top-level separators, that is, the number of tetrahedrons in  $(\cup_{i \in \mathcal{I}} \mathcal{H}_i) \cup (\cup_{i \notin \mathcal{I}} \mathcal{T}_i)$  which intersects a plane  $P_j$ , see Algorithm 15 line 4. For each  $i \in \mathcal{I} \stackrel{\text{def}}{=} \{i \in [k] : \alpha(\mathcal{T}_i) \leq n_i^{c_\alpha}\}$ , by Lemma 5.4.3, the number of tetrahedrons in  $\mathcal{H}_i$  intersect a plane  $P_j$  is at most

$$O\left(k^2 n_i^{2/3} \alpha_i^{-1/3} r_i^{-1/3}\right) = O\left(k^2 n_i^{2/3-c_r/3}\right),$$

since  $\alpha_i \geq 1$ . For each  $i \notin \mathcal{I}$ , by Lemma 5.1.4, the number of tetrahedrons in  $\mathcal{T}_i$  intersect a plane  $P_j$  is at most

$$O\left(k n_i^{2/3} \alpha_i^{-1/3}\right) = O\left(k n_i^{2/3-c_\alpha/3}\right)$$

since  $\alpha_i > n_i^{c_\alpha}$ . The two terms have same exponent for  $n$  when we set  $c_r = c_\alpha$ . Note Algorithm 13 requires that  $c_r + 2c_\alpha \leq 1$ . Thus, here we need  $c_r \leq 1/3$ .

Thus, the total number of tetrahedrons in  $(\cup_{i \in \mathcal{I}} \mathcal{H}_i) \cup (\cup_{i \notin \mathcal{I}} \mathcal{T}_i)$  which intersect a single separator plane  $P_j$  is then at most:

$$s \stackrel{\text{def}}{=} O\left(\sum_{1 \leq i \leq k} k^2 n_i^{2/3-c_r/3}\right) = O\left(k^{7/3+c_r/3} n^{2/3-c_r/3}\right).$$

The last inequality is by Jensen's inequality.

There are totally  $l$  top-level separators, which separates the whole truss into  $l + 1$  sep-

arate components and each component has  $O(n/l)$  vertices, according to Algorithm 15 line 3.

We plug these parameters into Lemma 2.5.2, the total fill-in size is

$$O\left(\left(\frac{n}{l}\right)^{4/3} l + s^2 l\right) = O(n^{4/3} l^{-1/3} + k^{14/3+2c_r/3} n^{4/3-2c_r/3} l),$$

and the multiplication count is

$$O\left(\left(\frac{n}{l}\right)^{2\omega/3} l + s^\omega l \log n\right) = O(n^{2\omega/3} l^{-2\omega/3+1} + k^{7\omega/3+\omega c_r/3} n^{2\omega/3-c_r\omega/3} l \log n).$$

This completes the proof. □

## 5.6 Running Times In Terms of Fast Matrix Multiplication

We now restate the running times of our algorithms in terms of faster matrix multiplication / inversion routines. Specifically, we assume inverting an  $n \times n$  matrix takes time  $O(n^\omega)$ , where  $\omega < 2.3728639$  [17].

We first examine purely nested dissection based algorithms. The running time of these algorithms are dominated by the cost of inverting the matrix at the top-most level. Thus, the running time of 3-D nested dissection from [49] as given in Theorem 2.5.1 is

$$O\left(n^{2\omega/3}\right),$$

while the performance of Lemma 2.5.2 becomes

$$O\left(n^{2\omega\beta/3+\alpha} + n^{\omega\gamma+\alpha}\right).$$

We now propagate these different costs for constructing the nested dissection partial states into our running time analyses.

For the bounded aspect ratio case described in Theorem 1.1.4, recall that the input truss is a union of  $k$  convex pieces of  $n_i$  vertices each. Following with parameter

$$r_i = n_i^{c_r}$$

now takes time

$$O\left(\sum_i \frac{n_i}{r_i} \cdot r_i^{2\omega/3}\right) = O\left(n^{1+c_r(2\omega/3-1)}\right),$$

while still giving a Schur complement of size

$$O\left(n^{1+c_r/3}\right)$$

on the boundaries, and total boundary size of

$$O\left(k^{c_r/3} n^{1-c_r/3}\right).$$

Theorem 2.5.1 then gives that solving this problem on just the boundary elements takes time

$$O\left(k^{2c_r\omega/9} n^{2(1-c_r/3)\omega/3}\right) = O\left(k^{2c_r\omega/9} n^{2\omega/3-2c_r\omega/9}\right),$$

and results in a total fill-in of

$$O\left(k^{4c_r/9} n^{4/3-4c_r/9}\right).$$

Putting these parameters back into the condition number bound of  $O(n^{2c_r})$  gives an iteration count of  $O(n^{c_r} \log(1/\epsilon))$ , which in turn gives a total cost of

$$\begin{aligned} & O\left(n^{1+c_r(2\omega/3-1)}\right) + O\left(k^{2c_r\omega/9} n^{2\omega/3-2c_r\omega/9}\right) + O\left(n^{c_r} \log(1/\epsilon)\right) \cdot [O\left(n^{1+c_r/3}\right) + O\left(k^{c_r/3} n^{4/3-4c_r/9}\right)] \\ &= O\left(n^{1+c_r(2\omega/3-1)} + k^{2c_r\omega/9} n^{2\omega/3-2c_r\omega/9} + n^{1+4c_r/3} \log(1/\epsilon) + k^{c_r/3} n^{4/3+5c_r/9} \log(1/\epsilon)\right). \end{aligned}$$

We can (slightly) simplify this using the fact that  $\omega \leq 3$  to drop the first term: it is always



upper bounded by the third. Also, since  $k \leq n$ , we will focus on optimizing the exponent on  $n$ , that is, we want to pick  $c_r$  to minimize the maximum of

$$\begin{aligned} & 2\omega/3 - 2c_r\omega/9 \\ & 1 + 4c_r/3 \\ & 4/3 + 5c_r/9 \end{aligned}$$

By running an LP solver, we get:

- when  $\omega = 3$  this is optimized at  $c_r = 1/2$ , which gives a total cost of  $O(k^{1/3}n^{5/3} \log(1/\epsilon))$ .
- when  $\omega = 2.3728639$ , this is optimized at  $c_r = 0.2295553$ . Here the exponents on  $n$  the three terms are 1.4608640, 1.3060737 and 1.4608640 respectively, and we have  $2\omega/9 > 1/3$ , so so the total cost is bounded by  $O(k^{0.1210452}n^{1.4608641} \log(1/\epsilon))$ .

For the more general case from Theorem 1.1.4, combining the bounds from Lemma 5.5.3 with the

- $O(n^{1+c_r(2\omega/3-1)})$  cost of computing the Schur complement of eliminating the innards of the hollowings, and
- the  $O(n^{1+c_r/3})$  size of the these Schur complements, and
- $O(n^{c_r} \log(1/\epsilon))$  iteration count of PCG

gives a total cost of:<sup>2</sup>

$$\begin{aligned} & O \left( n^{1+c_r(2\omega/3-1)} + n^{2\omega/3}l^{-2\omega/3+1} + k^{7\omega/3+c_r\omega/3}n^{2\omega/3-c_r\omega/3}l \right. \\ & \quad \left. + n^{c_r} \log(1/\epsilon) \left( n^{1+c_r/3} + n^{4/3}l^{-1/3} + k^{14/3+2c_r/3}n^{4/3-2c_r/3}l \right) \right). \end{aligned}$$

---

<sup>2</sup>We drop the  $\log n$  factor here, given  $\log n \ll n^c$  for any constant  $c > 0$ .

Since  $\omega \leq 3$ , we drop the first term. We can simplify this by moving  $k$  and  $\log(1/\epsilon)$  to the outermost, and only optimizing the remaining terms:

$$O\left(k^{7\omega/3+c_r\omega/3} \log(1/\epsilon)\right) \cdot \left(n^{2\omega/3} l^{-2\omega/3+1} + n^{2\omega/3-c_r\omega/3} l + \left(n^{1+4c_r/3} + n^{4/3+c_r} l^{-1/3} + n^{4/3+c_r/3} l\right)\right).$$

Let  $l = n^{c_l}$ . Since Algorithm 13 requires that  $c_r + 2c_\alpha \leq 1$  and in Lemma 5.5.3 we set  $c_r = c_\alpha$ , we have  $0 \leq c_r \leq 1/3$ . Subject to this requirement, we minimize the maximum of the following terms:

$$\begin{aligned} & 1 + \left(\frac{2\omega}{3} - 1\right) c_r \\ & \frac{2\omega}{3} + \left(-\frac{2\omega}{3} + 1\right) c_l, \\ & \frac{2\omega}{3} - \frac{\omega}{3} c_r + c_l, \\ & 1 + \frac{4}{3} c_r, \\ & \frac{4}{3} + c_r - \frac{1}{3} c_l, \\ & \frac{4}{3} + \frac{1}{3} c_r + c_l. \end{aligned}$$

By running an LP solver, we get:

- when  $\omega = 3$  this is optimized at  $c_r = 1/3$  and  $c_l = 1/6$ , which gives a total cost of  $O(k^{22/3} n^{11/6} \log(1/\epsilon))$  ( $11/6 \approx 1.8333$ )
- when  $\omega = 2.3728639$ , an optimum solution is  $c_r = 0.2210963$  and  $c_l = 0.1105482$  for a total cost of  $O(k^{5.7115596} n^{1.5175803} \log(1/\epsilon))$ .

## CHAPTER 6

### PARALLEL ALGORITHM FOR MIXED PACKING AND COVERING LPS

In this chapter, we prove Theorem 1.1.6.

#### 6.1 Converting an Optimization Problem to Feasibility Problems

Recall a mixed packing and covering LP has the form:

$$\min_{\mathbf{x} \geq \mathbf{0}} \{ \lambda : \mathbf{P}\mathbf{x} \leq \lambda \mathbf{p}, \mathbf{C}\mathbf{x} \geq \mathbf{c} \}.$$

By a standard binary search, one can reduce  $(1 + \epsilon)$ -approximately solving the above optimization problem into approximately solving  $\tilde{O}(1)$  feasibility problems (e.g., see [13]). Specifically, the feasibility problem is the following: either find  $\mathbf{x} \geq \mathbf{0}$  such that

$$0 < \max \mathbf{P}\mathbf{x} \leq (1 + \epsilon) \min \mathbf{C}\mathbf{x} \tag{6.1}$$

or conclude the following LP is infeasible

$$\begin{aligned} \mathbf{C}\mathbf{x} &\geq \mathbf{1} \\ \mathbf{P}\mathbf{x} &\leq (1 - 10\epsilon)\mathbf{1} \\ \mathbf{x} &\geq \mathbf{0} \end{aligned} \tag{6.2}$$

To certify that (6.2) is infeasible, we rely on its dual LP.

**Lemma 6.1.1.** *By duality, (6.2) is infeasible if there exists  $\mathbf{y}, \mathbf{z} \geq \mathbf{0}$  s.t.*

$$(1 - 10\epsilon) \frac{\mathbf{C}^\top \mathbf{z}}{\mathbf{1}^\top \mathbf{z}} < \frac{\mathbf{P}^\top \mathbf{y}}{\mathbf{1}^\top \mathbf{y}}. \tag{6.3}$$

*Proof.* Eqn. (6.3) is a direct reformulation of the dual LP of (6.2). Since we only need the sufficient condition, the result is by weak duality. If there exists any  $\mathbf{x} \geq 0$  satisfying (6.2), we have

$$\begin{aligned} (1 - 10\epsilon) \frac{\mathbf{x}^\top \mathbf{C}^\top \mathbf{z}}{\mathbf{1}^\top \mathbf{z}} &\geq (1 - 10\epsilon) \frac{\mathbf{1}^\top \mathbf{z}}{\mathbf{1}^\top \mathbf{z}} = 1 - 10\epsilon, \\ \frac{\mathbf{x}^\top \mathbf{P}^\top \mathbf{y}}{\mathbf{1}^\top \mathbf{y}} &\leq (1 - 10\epsilon) \frac{\mathbf{1}^\top \mathbf{y}}{\mathbf{1}^\top \mathbf{y}} = 1 - 10\epsilon. \end{aligned}$$

Together they give  $1 < 1$ , contradiction. □

## 6.2 Algorithm for the Feasibility Problem

Our algorithm follows the Lagrangian-relaxation approach as in [13, 78]. The idea is to convert all hard constraints into a single scalar-valued continuous and smoothed function. At each point, we can use a linear function to approximate this smoothed function in a small region.

### 6.2.1 Potential Function

As in [13], we use the soft-max  $\text{lmax}(\mathbf{P}\mathbf{x})$  and soft-min  $\text{lmin}(\mathbf{C}\mathbf{x})$  to write all constraints into one term. Recall that

$$\text{lmax}(\mathbf{P}\mathbf{x}) \stackrel{\text{def}}{=} \ln \left( \sum_j \exp(\mathbf{P}\mathbf{x})_j \right) = \ln \left( \sum_j \exp(\mathbf{P}_j^\top \mathbf{x}) \right),$$

and

$$\text{lmin}(\mathbf{C}\mathbf{x}) \stackrel{\text{def}}{=} -\ln \left( \sum_j \exp(-\mathbf{C}\mathbf{x})_j \right) = -\ln \left( \sum_j \exp(-\mathbf{C}_j^\top \mathbf{x}) \right),$$

where  $\mathbf{P}_j^\top, \mathbf{C}_j^\top$  are the  $j$ -th row of  $\mathbf{P}, \mathbf{C}$  respectively. In particular, these functions give smooth approximation to  $\max \mathbf{P}\mathbf{x}$  and  $\min \mathbf{C}\mathbf{x}$ :

$$\begin{aligned}\max \mathbf{P}\mathbf{x} &\leq \text{lmax}(\mathbf{P}\mathbf{x}) \leq \max \mathbf{P}\mathbf{x} + \ln n \\ \min \mathbf{C}\mathbf{x} &\geq \text{lmin}(\mathbf{C}\mathbf{x}) \geq \min \mathbf{C}\mathbf{x} - \ln n.\end{aligned}\tag{6.4}$$

When  $\text{lmax}(\mathbf{P}\mathbf{x}) \approx \text{lmin}(\mathbf{C}\mathbf{x})$ , we will have

$$\min \mathbf{C}\mathbf{x} \leq \max \mathbf{P}\mathbf{x} \leq \min \mathbf{C}\mathbf{x} + 2 \ln n.$$

If we can control that  $\min \mathbf{C}\mathbf{x} = \Theta(\ln n/\epsilon)$ , then the above inequality gives  $\max \mathbf{P}\mathbf{x} \leq (1 + \epsilon) \min \mathbf{C}\mathbf{x}$ .

### 6.2.2 Algorithm

We present our parallel  $\tilde{O}(1/\epsilon^3)$  routine in Algorithm 16 for solving the  $(1 + O(\epsilon))$ -feasibility problem above, that is, either find  $\mathbf{x} \geq 0$  satisfying (6.1), or certify the infeasibility of (6.2). The input contains a packing constraint matrix  $\mathbf{P} \in \mathbb{R}_{\geq 0}^{n_P \times m}$ , a covering constraint matrix  $\mathbf{C} \in \mathbb{R}_{\geq 0}^{n_C \times m}$ , and an error parameter  $0 < \epsilon$ . That is, there are  $m$  variables,  $n_P$  packing constraints, and  $n_C$  covering constraints. We also use  $n = n_P + n_C$  to denote the total number of constraints.

Algorithm 16 starts with small  $\mathbf{x}_i^{(0)} = \frac{1}{m\|\mathbf{P}_{:,i}\|_\infty}, \forall i \in [m]$ , and keeps increasing  $\mathbf{x}$  properly, until it reaches the terminate condition in line 4, that is,  $\max\{\max \mathbf{P}\mathbf{x}, \min \mathbf{C}\mathbf{x}\} \geq K = \frac{10 \ln n}{\epsilon}$ . The reason of the chosen  $K$  value is stated in Lemma 6.3.4. Roughly, when the difference of  $\min \mathbf{C}\mathbf{x}$  and  $\max \mathbf{P}\mathbf{x}$  becomes  $\epsilon$  factor smaller than  $\max\{\max \mathbf{P}\mathbf{x}, \min \mathbf{C}\mathbf{x}\}$ , we know that  $\max \mathbf{P}\mathbf{x} \leq (1 + O(\epsilon)) \min \mathbf{C}\mathbf{x}$ .

In each iteration of the while-loop, we first delete all covering constraints which has already reached  $K$ . Since  $\mathbf{x}$  never decreases, we know that once a row is deleted, we no longer need to look at it. Note the covering matrix cannot be empty, since we enter the

iteration with  $\min \mathbf{C}\mathbf{x} < K$ . We compute the vectors  $\mathbf{y}, \mathbf{z}$ , which are exponentials of the values of the packing and covering constraints respectively. We then compute  $\mathbf{a}$  and  $\mathbf{b}$ , which can be considered as gradients of  $\text{lmax}(\mathbf{P}\mathbf{x})$  and  $\text{lmin}(\mathbf{C}\mathbf{x})$  respectively, and use them to guide our update on  $\mathbf{x}$ . In particular, we update  $x_i$  if  $a_i \leq (1 - \epsilon/50)b_i$  (i.e.  $i \in B$ ). Furthermore, we update  $x_i$  multiplicatively by a factor depends on the ratio of  $\frac{a_i}{b_i}$ , as specified in Eqn. (6.5) and line 11. Note that the smallest update in our algorithm is by a factor of  $(1 + \Omega(\frac{\epsilon^2}{\ln n}))$ , which is the same as the fixed update step size in [78], and in general our updates take larger steps.

Note that in our analysis, we equivalently view  $\mathbf{z}$  as the full  $n_C$ -dimensional vector, where the coordinates corresponding to deleted constraints are filled by 0's. In particular, the matrix-vector product of the original  $\mathbf{C}$  with the  $n_C$ -dimensional  $\mathbf{z}$  will be the same as the product of the reduced covering matrix  $\mathbf{C}^{(t)}$  and reduced  $\mathbf{z}$ .

## 6.3 Analysis

### 6.3.1 Proof of Correctness

In this section we will show Algorithm 16 will terminate, and output the correct answer.

Lemma 6.3.1 shows that empty  $B$  certifies the infeasibility of the input instance (6.2), which proves the correctness if we end up in the case of line 9.

**Lemma 6.3.1.** *If the problem instance (6.2) is feasible, then*

$$\forall \mathbf{x} \geq 0, B = \{i : a_i \leq (1 - \frac{\epsilon}{50})b_i\} \neq \emptyset.$$

*Proof.* Let  $\mathbf{x}^*$  be a feasible solution of (6.2). Assume by contradiction,  $\exists \mathbf{x} \geq 0, \forall i \in [m], a_i^{(t)} > (1 - \frac{\epsilon}{50})b_i^{(t)}$ . By definition of  $\mathbf{a}, \mathbf{b}$ , it is equivalent to  $\exists \mathbf{y}, \mathbf{z} \geq 0$  such that

$$(1 - \frac{\epsilon}{50}) \frac{\mathbf{C}^\top \mathbf{z}}{\mathbf{1}^\top \mathbf{z}} < \frac{\mathbf{P}^\top \mathbf{y}}{\mathbf{1}^\top \mathbf{y}}.$$

---

**Algorithm 16** Parallel algorithm for mixed packing and covering LPs

---

**Input:**  $P, C, \epsilon$ **Output:** “infeasible” or  $\mathbf{x} \geq 0$  s.t.  $P\mathbf{x} \leq (1 + \epsilon)C\mathbf{x}$ 

- 1: Let  $K = \frac{10 \ln n}{\epsilon}$ ,  $\alpha = \frac{1}{K}$ , where  $n$  is the number of constraints.
  - 2: Initialize  $\mathbf{x}_i^{(0)} = \frac{1}{m \|P_{:,i}\|_\infty}$ ,  $\forall i \in [m]$ , where  $m$  is the number of variables.
  - 3: Let  $t = 0$ .
  - 4: **while**  $\max P\mathbf{x} < K$  and  $\min C\mathbf{x} < K$  **do**
  - 5:   Let  $C^{(t)}$  be  $C$  with rows  $j$  such that  $(C\mathbf{x}^{(t)})_j \geq K$  deleted.
  - 6:   Let  $\mathbf{y}^{(t)} = \exp(P\mathbf{x}^{(t)})$ ,  $\mathbf{z}^{(t)} = \exp(-C^{(t)}\mathbf{x}^{(t)})$ .
  - 7:    $\mathbf{a}^{(t)} = \frac{P^\top \mathbf{y}^{(t)}}{\mathbf{1}^\top \mathbf{y}^{(t)}}$ ,  $\mathbf{b}^{(t)} = \frac{(C^{(t)})^\top \mathbf{z}^{(t)}}{\mathbf{1}^\top \mathbf{z}^{(t)}}$ .
  - 8:   Define  $B^{(t)} = \{i : \mathbf{a}_i^{(t)} \leq (1 - \frac{\epsilon}{50})\mathbf{b}_i^{(t)}\}$ .
  - 9:   If  $B^{(t)} = \emptyset$ , then return “infeasible”.
  - 10:   Let
$$\Delta_i^{(t)} = \begin{cases} \frac{1}{2}(1 - \frac{\mathbf{a}_i^{(t)}}{\mathbf{b}_i^{(t)}}) \in [\epsilon/100, \frac{1}{2}] & \text{if } i \in B^{(t)} \\ 0 & \text{if } i \notin B^{(t)} \end{cases} \quad (6.5)$$
  - 11:    $\mathbf{x}_i^{(t+1)} \leftarrow \mathbf{x}_i^{(t)}(1 + \alpha \Delta_i^{(t)})$ .
  - 12:    $t \leftarrow t + 1$ .
  - 13: **end while**
  - 14: **return**  $\mathbf{x} = \frac{\mathbf{x}^{(t)}}{K}$ .
- 

Then the result follows directly from Lemma 6.1.1. □

If Algorithm 16 doesn't terminate with line 9, it must increase at least one variable by at least a factor of  $(1 + \frac{\epsilon^2}{10 \ln n})$  each iteration, so the algorithm must reach the termination condition of the while loop at some point, and we need to show the output  $\mathbf{x}$  satisfies (6.1).

We consider the following potential function,

$$f(\mathbf{x}^{(t)}) = \text{lmax}(P\mathbf{x}^{(t)}) - \text{lmin}(C^{(t)}\mathbf{x}^{(t)}) = \ln(\mathbf{1}^\top \mathbf{y}^{(t)}) + \ln(\mathbf{1}^\top \mathbf{z}^{(t)}).$$

We first quantify the changes of lmax and lmin when we update the variables. This type of smoothness analysis is standard in analyzing algorithms that make updates using gradient information. Similar results are derived in other works on packing and covering (See [54, 13]). The particular analysis we develop can deal with larger gradient steps. In particular,

the approach of our analysis allows updates that may move the gradients of some variables out of their respective coordinate-wise smooth regions, as long as we can still bound the combined impact on the potential function from updates of all variables. This approach can extend straightforwardly to show larger updates also work in [54], and improve their pure packing algorithm to run in  $\tilde{O}(\frac{N}{\epsilon^2})$  iterations.

**Lemma 6.3.2.** *At each iteration  $t$ ,*

$$\text{lmax}(\mathbf{P}\mathbf{x}^{(t+1)}) \leq \text{lmax}(\mathbf{P}\mathbf{x}^{(t)}) + \alpha \langle \mathbf{a}^{(t)}, (1 + \Delta^{(t)}) \cdot \Delta^{(t)} \cdot \mathbf{x}^{(t)} \rangle$$

and

$$\text{lmin}(\mathbf{C}^{(t+1)}\mathbf{x}^{(t+1)}) \geq \text{lmin}(\mathbf{C}^{(t)}\mathbf{x}^{(t)}) + \alpha \langle \mathbf{b}^{(t)}, (1 - \Delta^{(t)}) \cdot \Delta^{(t)} \cdot \mathbf{x}^{(t)} \rangle,$$

where  $\Delta \cdot \mathbf{x}$  is the entry-wise product vector, i.e.,  $(\Delta \cdot \mathbf{x})_i = \Delta_i x_i$ .

*Proof.* To simplify, we omit superscript  $(t)$  in the proof.

$$\begin{aligned} \text{lmax}(\mathbf{P}\mathbf{x}^{(t+1)}) &= \ln \sum_j \exp(\mathbf{P}_j^\top (\mathbf{x} + \alpha \Delta \cdot \mathbf{x})) \\ &= \ln \sum_j \exp(\mathbf{P}_j^\top \mathbf{x}) \cdot \exp(\alpha \mathbf{P}_j^\top (\Delta \cdot \mathbf{x})) \\ &\leq \ln \sum_j \exp(\mathbf{P}_j^\top \mathbf{x}) (1 + \alpha \mathbf{P}_j^\top (\Delta \cdot \mathbf{x}) + \alpha^2 (\mathbf{P}_j^\top (\Delta \cdot \mathbf{x}))^2). \end{aligned}$$

Recall  $\mathbf{P}_j^\top$  is the  $j$ -th row of  $\mathbf{P}$  (i.e. the  $j$ -th packing constraint). The last inequality is by Taylor expansion, with  $\epsilon \leq \Delta_i \leq \frac{1}{2}$ ,  $\alpha = \epsilon/10 \log n$ , and  $\mathbf{P}_j^\top \mathbf{x} \leq 10 \log n/\epsilon$  for all  $j$ , so  $\alpha \mathbf{P}_j^\top (\Delta \cdot \mathbf{x}) \leq \frac{1}{2}$ .



We can control the second order term as follows

$$\begin{aligned}
\alpha^2(\mathbf{P}_j^\top(\Delta \cdot \mathbf{x}))^2 &= \alpha^2 \left( \sum_i \mathbf{P}_{ji}(\Delta_i \mathbf{x}_i) \right)^2 \\
&\leq \alpha^2 \left( \sum_i \Delta_i \mathbf{P}_{ji}(\Delta_i \mathbf{x}_i) \right) \left( \sum_i \mathbf{P}_{ji} \mathbf{x}_i \right) \\
&\leq \alpha^2 \left( \sum_i \Delta_i \mathbf{P}_{ji}(\Delta_i \mathbf{x}_i) \right) \frac{10 \ln n}{\epsilon} \\
&= \alpha \left( \sum_i \Delta_i \mathbf{P}_{ji}(\Delta_i \mathbf{x}_i) \right) = \alpha \mathbf{P}_j^\top(\Delta \cdot \Delta \cdot \mathbf{x}).
\end{aligned}$$

The first inequality is by the Cauchy-Schwarz inequality  $\langle u, v \rangle^2 \leq \|u\|^2 \|v\|^2$ , with  $u_i = \sqrt{\Delta_i P_{ji}(\Delta_i x_i)}$  and  $v_i = \sqrt{P_{ji} x_i}$ . The second inequality is due to  $\mathbf{P}_j^\top \mathbf{x} \leq \frac{10 \ln n}{\epsilon}$  for all  $j$ . The last line is by our choice of  $\alpha = \frac{\epsilon}{10 \ln n}$ .

So far we have bounded the impact of the updates on each individual constraint, and we have

$$\text{lmax}(\mathbf{P}\mathbf{x}^{(t+1)}) \leq \ln \sum_j \exp(\mathbf{P}_j^\top \mathbf{x}) (1 + \alpha \mathbf{P}_j^\top(\Delta \cdot \mathbf{x}) + \alpha \mathbf{P}_j^\top(\Delta \cdot \Delta \cdot \mathbf{x})).$$

We then translate the changes on each constraint to the combined change on  $\text{lmax}(\mathbf{P}\mathbf{x})$ . Intuitively, the combined change is a convex combination on the changes of each constraint, weighted by their exponential values  $\mathbf{y}_j = \exp(\mathbf{P}_j^\top \mathbf{x})$ .

$$\begin{aligned}
&\ln \sum_j \exp(\mathbf{P}_j^\top \mathbf{x}) (1 + \alpha \mathbf{P}_j^\top(\Delta \cdot \mathbf{x}) + \alpha \mathbf{P}_j^\top(\Delta \cdot \Delta \cdot \mathbf{x})) \\
&= \ln \sum_j \mathbf{y}_j (1 + \alpha \mathbf{P}_j^\top(\Delta \cdot \mathbf{x}) + \alpha \mathbf{P}_j^\top(\Delta \cdot \Delta \cdot \mathbf{x})) \\
&= \ln \left( (\mathbf{1}^\top \mathbf{y}) \left( 1 + \sum_j \frac{\mathbf{y}_j}{\mathbf{1}^\top \mathbf{y}} (\alpha \mathbf{P}_j^\top(\Delta \cdot \mathbf{x}) + \alpha \mathbf{P}_j^\top(\Delta \cdot \Delta \cdot \mathbf{x})) \right) \right) \\
&= \ln \left( (\mathbf{1}^\top \mathbf{y}) \left( 1 + \alpha \left\langle \frac{\mathbf{y}}{\mathbf{1}^\top \mathbf{y}}, \mathbf{P}((\mathbf{1} + \Delta) \cdot \Delta \cdot \mathbf{x}) \right\rangle \right) \right).
\end{aligned}$$

We can then write out the change of  $\text{lmax}(\mathbf{P}\mathbf{x})$  explicitly as

$$\begin{aligned}
\text{lmax}(\mathbf{P}\mathbf{x}^{(t+1)}) &\leq \ln \left( (\mathbf{1}^\top \mathbf{y}) \left( 1 + \alpha \left\langle \frac{\mathbf{y}}{\mathbf{1}^\top \mathbf{y}}, \mathbf{P}((\mathbf{1} + \Delta) \cdot \Delta \cdot \mathbf{x}) \right\rangle \right) \right) \\
&= \ln \left( (\mathbf{1}^\top \mathbf{y}) \left( 1 + \alpha \left\langle \frac{\mathbf{P}^\top \mathbf{y}}{\mathbf{1}^\top \mathbf{y}}, (\mathbf{1} + \Delta) \cdot \Delta \cdot \mathbf{x} \right\rangle \right) \right) \\
&= \ln \left( (\mathbf{1}^\top \mathbf{y}) (1 + \alpha \langle \mathbf{a}, (\mathbf{1} + \Delta) \cdot \Delta \cdot \mathbf{x} \rangle) \right) \\
&= \text{lmax}(\mathbf{P}\mathbf{x}) + \ln (1 + \alpha \langle \mathbf{a}, (\mathbf{1} + \Delta) \cdot \Delta \cdot \mathbf{x} \rangle) \\
&\leq \text{lmax}(\mathbf{P}\mathbf{x}) + \alpha \langle \mathbf{a}, (\mathbf{1} + \Delta) \cdot \Delta \cdot \mathbf{x} \rangle.
\end{aligned}$$

Recall  $\mathbf{a} = \frac{\mathbf{P}^\top \mathbf{y}}{\mathbf{1}^\top \mathbf{y}}$  as defined in (7), and the last line is by  $\ln(1 + x) \leq x$  for  $x \geq 0$ .

For the  $\text{lmin}(\mathbf{C}^{(t)}\mathbf{x})$  part, we follow the same approach.

$$-\text{lmin}(\mathbf{C}^{(t+1)}\mathbf{x}^{(t+1)}) = \ln \sum_k \exp(-\mathbf{C}^{(t+1)}\mathbf{x}^{(t+1)})_k,$$

since  $\mathbf{C}^{(t+1)}$  can only have same or more rows dropped from  $\mathbf{C}^{(t)}$  due to the overly satisfied constraints, we know  $\ln \sum_k \exp(-\mathbf{C}^{(t+1)}\mathbf{x}^{(t+1)})_k \leq \ln \sum_k \exp(-\mathbf{C}^{(t)}\mathbf{x}^{(t+1)})_k$ . Again omitting the superscript  $(t)$ , we have

$$\begin{aligned}
-\text{lmin}(\mathbf{C}^{(t+1)}\mathbf{x}^{(t+1)}) &\leq \ln \sum_k \exp(-\mathbf{C}(\mathbf{x} + \alpha \Delta \cdot \mathbf{x}))_k \\
&= \ln \sum_k \exp(-\mathbf{C}\mathbf{x})_k \exp(-\alpha \mathbf{C}(\Delta \cdot \mathbf{x}))_k \\
&\leq \ln \sum_k \exp(-\mathbf{C}\mathbf{x})_k (1 - \alpha \mathbf{C}(\Delta \cdot \mathbf{x}) + \alpha^2 (\mathbf{C}(\Delta \cdot \mathbf{x}))^2)_k.
\end{aligned}$$

The last inequality is due to  $(\alpha \mathbf{C}(\Delta \cdot \mathbf{x}))_k \leq \frac{1}{2}$  for all  $k$ , since in  $\mathbf{C}^{(t)}$  we only keep those constraints that are not above  $\frac{10 \ln n}{\epsilon}$  yet, i.e.  $(\mathbf{C}\mathbf{x})_k \leq \frac{1}{\alpha}$  for all  $k$ . Again using Cauchy-Schwarz inequality as before and derivations similar to the  $\text{lmax}(\mathbf{P}\mathbf{x})$  case, we get

$$-\text{lmin}(\mathbf{C}^{(t+1)}\mathbf{x}^{(t+1)}) \leq -\text{lmin}(\mathbf{C}^{(t)}\mathbf{x}^{(t)}) - \alpha \langle \mathbf{b}, (1 - \Delta) \cdot \Delta \cdot \mathbf{x} \rangle.$$

□

With the above bounds on the changes of the two components  $\text{lmax}(\mathbf{P}\mathbf{x})$  and  $\text{lmin}(\mathbf{C}\mathbf{x})$ , we can show how our updates move the potential function  $f(\mathbf{x})$ .

**Lemma 6.3.3.** *Given  $\max \mathbf{P}\mathbf{x}^{(t)} < \frac{10 \ln n}{\epsilon}$  and  $\min \mathbf{C}\mathbf{x}^{(t)} < \frac{10 \ln n}{\epsilon}$ , we always have  $f(\mathbf{x}^{(t)}) \leq 2 \ln n$  during the execution of Algorithm 16.*

*Proof.* Initially,  $\mathbf{x}_i^{(0)} = \frac{1}{m\|\mathbf{P}_{:,i}\|_\infty}$ , we have  $\mathbf{P}\mathbf{x}^{(0)} \leq \mathbf{1}$  and  $\mathbf{C}\mathbf{x} \geq 0$ , thus  $f(\mathbf{x}^{(0)}) \leq 2 \ln n$ . To show  $f(\mathbf{x}) \leq 2 \ln n$  for all iterations  $t$  before terminate, it suffices to show that  $f(\mathbf{x})$  is non-increasing during the process. From Lemma 6.3.2,

$$\begin{aligned} f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)}) &\leq \alpha \langle \mathbf{a}, (1 + \Delta) \cdot \Delta \cdot \mathbf{x}^{(t)} \rangle - \alpha \langle \mathbf{b}, (1 - \Delta) \cdot \Delta \cdot \mathbf{x}^{(t)} \rangle \\ &= \sum_i \alpha \Delta_i \mathbf{x}_i (\mathbf{a}_i (1 + \Delta_i) - \mathbf{b}_i (1 - \Delta_i)). \end{aligned}$$

For each  $i \in [m]$ , by our update rule (6.5), either  $\Delta_i$ , or  $\Delta_i = \frac{1}{2}(1 - \frac{\mathbf{a}_i}{\mathbf{b}_i})$ , in which case

$$\mathbf{a}_i(1 + \Delta_i) - \mathbf{b}_i(1 - \Delta_i) = \frac{3\mathbf{a}_i\mathbf{b}_i - \mathbf{a}_i^2}{2\mathbf{b}_i} - \frac{\mathbf{a}_i + \mathbf{b}_i}{2} = \frac{2\mathbf{a}_i\mathbf{b}_i - \mathbf{a}_i^2 - \mathbf{b}_i^2}{2\mathbf{b}_i} \leq 0,$$

so all the summands are non-positive, thus  $f(\mathbf{x})$  is non-increasing. □

The above lemma guarantees that the difference between  $\text{lmax}(\mathbf{P}\mathbf{x})$  and  $\text{lmin}(\mathbf{C}\mathbf{x})$  is bounded by  $2 \ln n$ , which by Eqn. (6.4) suggests  $\max \mathbf{P}\mathbf{x} \leq \min \mathbf{C}\mathbf{x} + O(\ln n)$ . Then when the two terms are large at termination, we are approximately feasible as the difference is a factor of  $\epsilon$  smaller.

**Lemma 6.3.4.** *If Algorithm 16 terminates with line 14, then it returns an  $\mathbf{x} \geq 0$  with  $0 < \max \mathbf{P}\mathbf{x} \leq (1 + \epsilon) \min \mathbf{C}\mathbf{x}$ .*

*Proof.* Suppose the algorithm terminates at iteration  $T$ , that is,  $\max \mathbf{P}\mathbf{x}^{(T)} \geq \frac{10 \ln n}{\epsilon}$  or  $\min \mathbf{C}\mathbf{x}^{(T)} \geq \frac{10 \ln n}{\epsilon}$ . Consider iteration  $T - 1$ , the covering matrix is not empty (otherwise,

the algorithm terminates before iteration  $T$ ). Since  $\mathbf{x}^{(T)} = \mathbf{x}^{(T-1)} \cdot (1 + \alpha\Delta^{(T-1)}) \leq (1 + \frac{5\epsilon}{\ln n})\mathbf{x}^{(T-1)}$ , we have  $\max \mathbf{P}\mathbf{x}^{(T-1)} \geq \frac{5\ln n}{\epsilon}$  or  $\min \mathbf{C}\mathbf{x}^{(T-1)} \geq \frac{5\ln n}{\epsilon}$ .

By Lemma 6.3.3,

$$\max \mathbf{P}\mathbf{x}^{(T-1)} \leq \text{lmax}(\mathbf{P}\mathbf{x}^{(T-1)}) \leq \text{lmin}(\mathbf{C}^{(T-1)}\mathbf{x}^{(T-1)}) + 2\ln n \leq \min \mathbf{C}\mathbf{x}^{(T-1)} + 2\ln n.$$

Since  $2\ln n \leq \epsilon \cdot \frac{5\ln n}{\epsilon}$ , we have

$$\max \mathbf{P}\mathbf{x}^{(T-1)} \leq (1 + \epsilon) \min \mathbf{C}\mathbf{x}^{(T-1)}.$$

This also gives  $\max \mathbf{P}\mathbf{x}^{(T)} \leq (1 + \epsilon) \min \mathbf{C}\mathbf{x}^{(T)}$ , since  $\mathbf{x}^{(T)}$  is within in multiplicative factor  $1 + \frac{\epsilon}{10\ln n}$  of  $\mathbf{x}^{(T-1)}$ . Since we start with  $\mathbf{x} > 0$ , and only increase  $\mathbf{x}$ , we also have  $\max \mathbf{P}\mathbf{x} > 0$ . So the  $\mathbf{x}$  we return at the end satisfies (6.1).  $\square$

### 6.3.2 Analysis of Convergence

So far we have proved that Algorithm 16 will terminate, and will either output  $\mathbf{x}$  satisfying (6.1) at the end, or terminate earlier and correctly certify (6.2) is infeasible. In this section we show that if (6.2) is feasible, Algorithm 16 must finish with the first case in  $\tilde{O}(\frac{1}{\epsilon^3})$  iterations, so if the algorithm takes more than  $\frac{1000\ln n \ln(\frac{m}{\epsilon})}{\epsilon^3}$  iterations to complete, we can terminate it, and correctly output that (6.2) is infeasible.

We adapt the concept *phase* from Young's algorithm. Note phase is only used in our analysis, and our algorithm does not contain phase. Formally, phase  $s$  contains the iterations  $t$  such that

$$\frac{n_P}{n_C} \cdot 2^s \leq \frac{\mathbf{1}^T \mathbf{y}^{(t)}}{\mathbf{1}^T \mathbf{z}^{(t)}} < \frac{n_P}{n_C} \cdot 2^{s+1}$$

where  $n_P$  is the number of packing constraints and  $n_C$  is the number of covering constraints.

Since we only increase  $\mathbf{x}$ ,  $\frac{\mathbf{1}^T \mathbf{y}}{\mathbf{1}^T \mathbf{z}}$  is monotonically increasing, so each phase covers a con-

secutive sequence of iterations. Furthermore, as  $\ln(\frac{1^T \mathbf{y}}{1^T \mathbf{z}}) = \text{lmax}(\mathbf{P}\mathbf{x}) + \text{lmin}(\mathbf{C}\mathbf{x})$  measures global progress towards termination, each phase captures a fixed amount of progress. From our definition of phases, and the termination condition, we have

**Lemma 6.3.5.** *The total number of phases in Algorithm 16 is  $O(\frac{\log n}{\epsilon})$ .*

*Proof.* Since  $\mathbf{x}$  is monotonically increasing,  $\mathbf{y} = \exp(\mathbf{P}\mathbf{x})$  and  $\mathbf{z} = \exp(-\mathbf{C}\mathbf{x})$  are monotonically increasing and decreasing respectively, which implies that the quantity  $\frac{1^T \mathbf{y}}{1^T \mathbf{z}}$  is monotonically increasing. Initially  $\mathbf{P}\mathbf{x}^{(0)} \geq 0$ ,  $\mathbf{C}\mathbf{x}^{(0)} \geq 0$ , we know  $\frac{1^T \mathbf{y}^{(0)}}{1^T \mathbf{z}^{(0)}} \geq \frac{n_P}{n_C}$ . By the termination condition in Algorithm 16, the ratio never goes beyond  $n_P \exp(\frac{10 \log n}{\epsilon})$ . Therefore, the total number of phases is  $O(\frac{\log n}{\epsilon})$ .  $\square$

We now bound the number of iterations in a single phase. The iterations of a phase are divided into two groups, the bad iterations and the good iterations, formally defined as follows.

**Definition 6.3.6.** *If in an iteration  $t$ , we have for all  $i$*

$$\frac{\mathbf{a}_i^{(t)}}{\mathbf{b}_i^{(t)}} > \frac{1}{3}, \quad (6.6)$$

*then we call it a good iteration. Otherwise we call it a bad iteration.*

Note a phase may contain only bad iterations or only good iterations. We bound the total number of iterations in the two groups separately.

As discussed earlier, the bad iterations capture the initial warm-up iterations of a phase, where in any bad iteration, we can identify some variable  $\mathbf{x}_i$  with a strong signal (i.e.  $\frac{\mathbf{a}_i}{\mathbf{b}_i} \leq \frac{1}{3}$ ), so we can increase the variable by a lot. This restricts the warm-up sequence from getting too long, and we formalize the intuition in the following lemma.

**Lemma 6.3.7.** *In a single phase, the number of bad iterations is at most  $O(\ln n \ln(\frac{m}{\epsilon})/\epsilon)$ .*

*Proof.* We will prove the result by showing that there cannot be any bad iteration after the initial  $100 \ln n \ln(\frac{m}{\epsilon})/\epsilon$  iterations of a phase. By contradiction, if for any variable  $i$ , after  $\Omega(\ln n \ln(\frac{m}{\epsilon})/\epsilon)$  iterations of a phase, we have at iteration  $t$  such that for some  $i$ ,

$$\frac{\mathbf{a}_i^{(t)}}{\mathbf{b}_i^{(t)}} = \frac{(\mathbf{P}^T \mathbf{y}^{(t)})_i}{\mathbf{1}^T \mathbf{y}^{(t)}} \frac{\mathbf{1}^T \mathbf{z}^{(t)}}{(\mathbf{C}^T \mathbf{z}^{(t)})_i} \leq \frac{1}{3},$$

then this ratio is at most  $\frac{2}{3}$  in all previous iterations of this phase, since  $\frac{(\mathbf{P}^T \mathbf{y})_i}{(\mathbf{C}^T \mathbf{z})_i}$  is monotonically increasing, and  $2^s \leq \frac{\mathbf{1}^T \mathbf{y}}{\mathbf{1}^T \mathbf{z}} < 2^{s+1}$  in this phase. Equivalently, this is saying  $\mathbf{a}_i \leq \frac{2}{3} \mathbf{b}_i$ , so  $i \in B$  in all previous  $\Omega(\ln n \ln(\frac{m}{\epsilon})/\epsilon)$  iterations of the phase, and  $\Delta_i \geq \frac{1}{6}$  in all those iterations.

Each iteration the multiplicative update on  $\mathbf{x}_i$  is  $(1 + \alpha \Delta_i)$ , which is  $(1 + \Theta(\frac{\epsilon}{10 \ln n}))$  since  $\Delta_i \geq \frac{1}{6}$ . As  $\mathbf{x}_i$  starts with  $\frac{1}{m \|P_{:,i}\|_\infty}$ , after  $100 \ln n \ln(\frac{m}{\epsilon})/\epsilon$  updates, we have  $\mathbf{x}_i >> \frac{10 \ln n}{\epsilon \|P_{:,i}\|_\infty}$ , which gives  $\max \mathbf{P} \mathbf{x} >> \frac{10 \ln n}{\epsilon}$ , so the algorithm must have terminated.  $\square$

The above lemma guarantees that all iterations after the first  $100 \ln n \ln(\frac{m}{\epsilon})/\epsilon$  must be good iterations, so we proceed to bound the number of these good iterations in a single phase. Without loss of generality, we index these good iterations in a phase as  $1, \dots, T$  by shifting  $t$ .

We first identify one variable that must be updated extensively in these iterations.

**Claim 6.3.8.** *Suppose the instance (6.2) is feasible, then There exists  $i \in [m]$  such that*

$$\sum_{t=1}^T \mathbf{b}_i^{(t)} - \mathbf{a}_i^{(t)} \geq 10\epsilon \sum_{t=1}^T \mathbf{b}_i^{(t)}. \quad (6.7)$$

*Proof.* Define  $\bar{\mathbf{y}}$  and  $\bar{\mathbf{z}}$  to be the sum of the normalized gradients of iterations  $1, \dots, T$ , that is,

$$\bar{\mathbf{y}} = \sum_{t=1}^T \frac{\mathbf{y}^{(t)}}{\mathbf{1}^T \mathbf{y}^{(t)}}, \quad \bar{\mathbf{z}} = \sum_{t=1}^T \frac{\mathbf{z}^{(t)}}{\mathbf{1}^T \mathbf{z}^{(t)}}.$$

Note  $\mathbf{1}^T \bar{\mathbf{y}} = \mathbf{1}^T \bar{\mathbf{z}} = T$ . Recall  $\mathbf{a}_i^{(t)}$  and  $\mathbf{b}_i^{(t)}$  are respectively  $\frac{(\mathbf{P}^T \mathbf{y}^{(t)})_i}{\mathbf{1}^T \mathbf{y}^{(t)}}$  and  $\frac{(\mathbf{C}^T \mathbf{z}^{(t)})_i}{\mathbf{1}^T \mathbf{z}^{(t)}}$ , then

$$\sum_{t=1}^T \mathbf{a}_i^{(t)} = \frac{T(\mathbf{P}^T \bar{\mathbf{y}})_i}{\mathbf{1}^T \bar{\mathbf{y}}}, \quad \sum_{t=1}^T \mathbf{b}_i^{(t)} = \frac{T(\mathbf{C}^T \bar{\mathbf{z}})_i}{\mathbf{1}^T \bar{\mathbf{z}}}.$$

Assume by contradiction,  $\forall i \in [m], \sum_{t=1}^T \mathbf{a}_i^{(t)} > (1 - 10\epsilon) \sum_{t=1}^T \mathbf{b}_i^{(t)}$ , that is,

$$\frac{\mathbf{P}^T \bar{\mathbf{y}}}{\mathbf{1}^T \bar{\mathbf{y}}} > (1 - 10\epsilon) \frac{\mathbf{C}^T \bar{\mathbf{z}}}{\mathbf{1}^T \bar{\mathbf{z}}}.$$

By Lemma 6.1.1,  $\bar{\mathbf{y}}, \bar{\mathbf{z}}$  certify infeasibility of the instance (6.2), which contradicts the assumption.  $\square$

The above claim gives us a variable that on average has smaller packing gradients than covering gradients in this iteration. Together with the property we have on the good iterations (6.6), we can bound the number of good iterations.

**Lemma 6.3.9.** *In a single phase, the number of good iterations is at most  $O(\ln n \ln(\frac{m}{\epsilon})/\epsilon^2)$ .*

*Proof.* Let  $\mathbf{x}_i$  be a variable satisfying Eqn. (6.7). We want to turn Eqn. (6.7) into some lower bound on the total multiplicative update on  $\mathbf{x}_i$  through these iterations. Intuitively, a bad case is that in some iteration  $t$ ,  $\mathbf{a}_i^{(t)}, \mathbf{b}_i^{(t)}$  are much larger than the values in other iterations, since they can dominate the terms from other iterations in Eqn. (6.7), but not much to the total update of  $\mathbf{x}_i$ , since their ratio is what matters to the update. However, since we are inside one single phase, and only looking at good iterations, we can show the bad scenario will not show up.

Formally, let  $l = \mathbf{a}_i^{(1)}$  and  $u = \mathbf{b}_i^{(1)}$ . Since  $(\mathbf{P}^T \mathbf{y})_i$  monotonically increases, and  $\mathbf{1}^T \mathbf{y}$  will increase but not by more than a factor of 2 in a phase, we have

$$\mathbf{a}_i^{(t)} = \frac{(\mathbf{P}^T \mathbf{y}^{(t)})_i}{\mathbf{1}^T \mathbf{y}^{(t)}} \geq l/2 \quad \forall t = 1, \dots, T \quad (6.8)$$

Similarly, we have

$$\mathbf{b}_i^{(t)} = \frac{(\mathbf{C}^T \mathbf{z}^{(t)})_i}{\mathbf{1}^T \mathbf{z}^{(t)}} \leq 2u \quad \forall t = 1, \dots, T \quad (6.9)$$

Furthermore, since we are looking at the good iterations, we have

$$l \geq \frac{1}{3}u.$$

The inequalities above allow us to turn the difference-based guarantee from Eqn. (6.7) into lower bounds on ratios we need.

By the update (6.5), we have

$$\Delta_i^{(t)} \geq \frac{(1 - \frac{\epsilon}{50})\mathbf{b}_i^{(t)} - \mathbf{a}_i^{(t)}}{2\mathbf{b}_i^{(t)}}.$$

So we can lower bound the total update on  $\mathbf{x}_i$  as follows

$$\begin{aligned} \mathbf{x}_i^{(T)} &\geq \mathbf{x}_i^{(1)} \exp \left( \frac{\alpha \sum_t \Delta_i^{(t)}}{2} \right) \\ &= \mathbf{x}_i^{(1)} \exp \left( \frac{\alpha}{4} \sum_t \frac{(1 - \frac{\epsilon}{50})\mathbf{b}_i^{(t)} - \mathbf{a}_i^{(t)}}{\mathbf{b}_i^{(t)}} \right) \\ &\geq \mathbf{x}_i^{(1)} \exp \left( \frac{\alpha}{4} \sum_t \left( \frac{\mathbf{b}_i^{(t)} - \mathbf{a}_i^{(t)}}{2u} - \frac{\epsilon}{50} \right) \right) \end{aligned}$$

where we used (6.9) in the last line.

From Eqn. (6.7), we have

$$\begin{aligned} \sum_t \mathbf{b}_i^{(t)} - \mathbf{a}_i^{(t)} &\geq 10\epsilon \sum_i \mathbf{b}_i^{(t)} \\ &\geq \frac{10\epsilon}{1 - 10\epsilon} \sum_t \mathbf{a}_i^{(t)} \\ &\geq \frac{\epsilon u T}{1 - 10\epsilon} \geq \epsilon u T. \end{aligned}$$



The first two lines both follow from Eqn. (6.7), the next line follows from  $\mathbf{a}_i^{(t)} \geq l/2 \geq u/6$ .

Thus

$$\mathbf{x}_i^{(T)} \geq \mathbf{x}_i^{(1)} \exp\left(\frac{\epsilon\alpha T}{8} - \frac{\epsilon\alpha T}{200}\right) \geq \frac{1}{m \|\mathbf{P}_{:,i}\|_\infty} \exp\left(\frac{\epsilon\alpha T}{10}\right).$$

If  $T \geq \frac{100 \ln n \ln \frac{m}{\epsilon}}{\epsilon^2} \geq \frac{100 \ln \frac{m}{\epsilon}}{\epsilon\alpha}$ , we have  $\mathbf{x}_i^{(T)} \gg \frac{10 \ln n}{\epsilon \|\mathbf{P}_{:,i}\|_\infty}$ . So the algorithm must have terminated since  $\max \mathbf{P}\mathbf{x} \gg \frac{10 \ln n}{\epsilon}$ .  $\square$

Lemma 6.3.7 and Lemma 6.3.9 bound the total number of iterations in a phase by  $\tilde{O}(\frac{1}{\epsilon^2})$ , together with the bound on the number of phases, which is  $\tilde{O}(\frac{1}{\epsilon})$ , we guarantee the total number of iterations in Algorithm 16 is  $\tilde{O}(\frac{1}{\epsilon^3})$  if the LP in (6.2) is feasible.

*Proof of Theorem 1.1.6.* The correctness and convergence follows from the lemmas in the prior sections. We only need to look at the running time and total work.

At each iteration, we compute all updated values in  $O(\log N)$  parallel time. Since the total number of iterations is  $\tilde{O}(\frac{1}{\epsilon^3})$ , the algorithm terminates in parallel time  $\tilde{O}(\frac{1}{\epsilon^3})$ .

To see the total work, consider the following implementation. For each  $i \in [m]$ , we maintain  $\mathbf{P}_{ji}\mathbf{x}_i$  if  $\mathbf{P}_{ij} \neq 0$ ; similarly we maintain  $\mathbf{C}_{ji}\mathbf{x}_i$  if  $\mathbf{C}_{ij} \neq 0$ . Besides, we maintain the values of  $\mathbf{y}, \mathbf{z}, \mathbf{P}^T \mathbf{y}, \mathbf{C}^T \mathbf{z}, \mathbf{1}^T \mathbf{y}$  and  $\mathbf{1}^T \mathbf{z}$ . When we update  $\mathbf{x}_i$ , we update these values accordingly, with work proportional to the number of non-zero entries in the  $i$ th column of the constraint matrix. For each fixed variable  $\mathbf{x}_i$ , the total time of updates is at most  $\tilde{O}(\frac{1}{\epsilon^2})$ . Thus, the work on this part is  $\tilde{O}(\frac{N}{\epsilon^2})$ . Additionally, we need to compute the  $\mathbf{a}_i, \mathbf{b}_i$  for all variables at the beginning of each iteration to determine which variables to update, this takes  $\tilde{O}(N)$  work each iteration, so the total work is  $\tilde{O}(\frac{N}{\epsilon^3})$ .

We see the majority of the work is actually on computing the gradients for the variables we may not update. We point out that we can implement the same lazy update as in [78], which on a high level is just that if a variable has a large  $\frac{a_i}{b_i}$  in an iteration, and is not updated, we don't recompute its gradients, until  $\frac{\mathbf{1}^T \mathbf{y}}{\mathbf{1}^T \mathbf{z}}$  grows by more than a factor of  $1 + \epsilon$ . This can reduce the work to  $\tilde{O}(\frac{N}{\epsilon^2})$ , but requires a centralized step to control the phases. We omit the details as it is a straightforward adaptation.  $\square$

### 6.3.3 Pure Packing and Pure Covering LPs

We point out that in the case of pure packing or pure covering LPs, Algorithm 16 converges in  $\tilde{O}(\frac{1}{\epsilon^2})$  iterations. This improves upon the result of [61], since our algorithm is deterministic, and does not need centralized steps.

Given packing LP in the optimization form

$$\max_{x \geq 0} \{ \mathbf{1}^T \mathbf{x} : \mathbf{P} \mathbf{x} \leq \mathbf{1} \}$$

via standard reduction and scaling, we need to solve a  $(1 + \epsilon)$ -feasibility problem the same as in the mixed packing and covering case specified in (6.1) and (6.2). In the case of pure packing, we will have  $\mathbf{C} = c\mathbf{1}^T$  for some constant  $c$ .

The correctness proof follows from the mixed case, and we discuss how we get faster convergence for pure packing.

The special structure of  $\mathbf{C}$  greatly simplifies the convergence analysis, as now  $z$  is a scalar, and  $\mathbf{b}_i^{(t)} = c$  for all variables  $\mathbf{x}_i$  across all iterations  $t$ . This allows us to aggregate the good iterations from all phases, and bound the total number of good phases by  $\tilde{O}(\frac{1}{\epsilon^2})$ , which will lead to  $\tilde{O}(\frac{1}{\epsilon^2})$  total iterations.

In particular, now we look at all the good iterations across all phases together, and WLOG number them  $1, \dots, T$ . Claim 6.3.8 still holds, as it does not rely on phases. Then we can prove a stronger version of Lemma 6.3.9

**Lemma 6.3.10.** *The total number of good iterations  $T$  is at most  $O(\ln n \ln(\frac{m}{\epsilon})/\epsilon^2)$ .*

The proof is a straightforward adaptation of the proof of Lemma 6.3.9. The proof is simpler, since now  $\mathbf{b}_i$  is a constant across all iterations, so the property (6.6) we have on the good iterations directly put all the values on the same scale, so we can lower bound the total update on  $\mathbf{x}_i$  across all good iterations.

# **Appendices**

## APPENDIX A

### PROOFS OF LINEAR ALGEBRA FACTS

#### A.1 Schur Complements

*Proof of Fact 2.1.2.* We expand the left hand side,

$$\begin{pmatrix} \mathbf{x}^\top & \mathbf{y}^\top \end{pmatrix} \mathbf{A} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \mathbf{x}^\top \mathbf{A}_{SS} \mathbf{x} + 2\mathbf{x}^\top \mathbf{A}_{ST} \mathbf{y} + \mathbf{y}^\top \mathbf{A}_{TT} \mathbf{y}. \quad (\text{A.1})$$

Taking derivative w.r.t.  $\mathbf{y}$  and setting it to be 0 give that

$$2\mathbf{A}_{TT} \mathbf{y} + 2\mathbf{A}_{ST}^\top \mathbf{x} = \mathbf{0}.$$

Plugging  $\mathbf{y} = -\mathbf{A}_{TT}^\dagger \mathbf{A}_{ST}^\top \mathbf{x}$  into (A.1),

$$\min_{\mathbf{y}} \begin{pmatrix} \mathbf{x}^\top & \mathbf{y}^\top \end{pmatrix} \mathbf{A} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \mathbf{x}^\top \mathbf{A}_{SS} \mathbf{x} - \mathbf{x}^\top \mathbf{A}_{ST} \mathbf{A}_{TT}^\dagger \mathbf{A}_{ST}^\top \mathbf{x} = \mathbf{x}^\top \text{SC}[\mathbf{A}]_T \mathbf{x}.$$

This completes the proof. □

To prove Fact 2.1.3, we need the following special case of Weyl inequalities.

**Theorem A.1.1** (A special case of Weyl inequalities). *Let  $\mathbf{H} \in \mathbb{R}^n$  and  $\mathbf{H} = \mathbf{H}_1 + \mathbf{H}_2$  where  $\mathbf{H}_1, \mathbf{H}_2$  are symmetric matrices and  $\mathbf{H}_2$  is a PSD matrix. Let  $\lambda_1(\cdot) \geq \dots \geq \lambda_n(\cdot)$  be eigenvalues of a matrix. Then,  $\lambda_i(\mathbf{H}) \geq \lambda_i(\mathbf{H}_1), i = 1, 2, \dots, n$ .*

*Proof of Fact 2.1.3.* The first statement immediately follows Fact 2.1.2.

To prove the second statement, we decompose  $\mathbf{A}$ :

$$\begin{aligned}\mathbf{A} &= \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{SC}[\mathbf{A}]_T \end{pmatrix} + \begin{pmatrix} \mathbf{A}_{SS} & \mathbf{A}_{ST} \\ \mathbf{A}_{ST}^\top & \mathbf{A}_{ST}^\top \mathbf{A}_{SS}^{-1} \mathbf{A}_{ST} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{SC}[\mathbf{A}]_T \end{pmatrix} + \begin{pmatrix} \mathbf{A}_{SS}^{1/2} \\ \mathbf{A}_{ST}^\top \mathbf{A}_{SS}^{-1/2} \end{pmatrix}^\top \begin{pmatrix} \mathbf{A}_{SS}^{1/2} \\ \mathbf{A}_{ST}^\top \mathbf{A}_{SS}^{-1/2} \end{pmatrix}.\end{aligned}$$

Here we assume that  $\mathbf{A}_{SS}$  is invertible, otherwise we can use pseudo-inverse. The first matrix is symmetric, and the second matrix is symmetric and PSD. By Theorem A.1.1,  $\lambda_{\max}(\mathbf{A}) \geq \lambda_{\max}(\text{SC}[\mathbf{A}]_T)$ .  $\square$

## A.2 Solutions of Linear Systems

In this section, we prove some facts on linear system solutions.

*Proof of Fact 2.2.4.* First observe

$$\begin{aligned}\|\mathbf{A}^\top \mathbf{A} \mathbf{x} - \mathbf{A}^\top \mathbf{c}\|_{(\mathbf{A}^\top \mathbf{A})^\dagger}^2 &= (\mathbf{A}^\top \mathbf{A} \mathbf{x} - \mathbf{A}^\top \mathbf{c})^\top (\mathbf{A}^\top \mathbf{A})^\dagger (\mathbf{A}^\top \mathbf{A} \mathbf{x} - \mathbf{A}^\top \mathbf{c}) \\ &= (\mathbf{A} \mathbf{x} - \mathbf{c})^\top \mathbf{\Pi}_A (\mathbf{A} \mathbf{x} - \mathbf{c}) \\ &= (\mathbf{\Pi}_A \mathbf{A} \mathbf{x} - \mathbf{\Pi}_A \mathbf{c})^\top (\mathbf{\Pi}_A \mathbf{A} \mathbf{x} - \mathbf{\Pi}_A \mathbf{c}) \\ &= \|\mathbf{A} \mathbf{x} - \mathbf{\Pi}_A \mathbf{c}\|_2^2.\end{aligned}$$

Taking  $\mathbf{x} = \mathbf{0}$  it follows that  $\|\mathbf{\Pi}_A \mathbf{c}\| = \|\mathbf{A}^\top \mathbf{c}\|_{(\mathbf{A}^\top \mathbf{A})^\dagger}$ , and combining our observations immediately gives the second part of the Fact.

A similar argument gives claims for  $\|\mathbf{x} - \mathbf{x}^*\|_{\mathbf{A}^\top \mathbf{A}}$ .

$$\begin{aligned}
\|\mathbf{x} - \mathbf{x}^*\|_{\mathbf{A}^\top \mathbf{A}}^2 &= (\mathbf{x} - \mathbf{x}^*)^\top \mathbf{A}^\top \mathbf{A} (\mathbf{x} - \mathbf{x}^*) \\
&= \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} - 2\mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x}^* + (\mathbf{x}^*)^\top \mathbf{A}^\top \mathbf{A} \mathbf{x}^* \\
&= \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} - 2\mathbf{x}^\top \mathbf{A}^\top \Pi_{\mathbf{A}} \mathbf{c} + \mathbf{c}^\top \Pi_{\mathbf{A}}^\top \Pi_{\mathbf{A}} \mathbf{c} \\
&= \|\mathbf{A} \mathbf{x} - \Pi_{\mathbf{A}} \mathbf{c}\|_2^2.
\end{aligned}$$

The third equality uses the fact that  $\mathbf{A} \mathbf{x}^* = \Pi_{\mathbf{A}} \mathbf{c}$ . This completes the proof.  $\square$

*Proof of Lemma 2.2.5.* As  $\|\mathbf{A} \mathbf{c}\|_2^2$  is integral, the condition of  $\mathbf{A}^\top \mathbf{c} \neq 0$  implies  $\|\mathbf{A}^\top \mathbf{c}\|_2 \geq 1$ .

1. Consider the SVD of  $\mathbf{A}$ :

$$\mathbf{A} = \sum_i \sigma_i \mathbf{r}^i (\mathbf{s}^i)^\top.$$

Substituting it in for  $\mathbf{A}^\top$  gives

$$\mathbf{A}^\top \mathbf{c} = \sum_i (\sigma_i (\mathbf{r}^i)^\top \mathbf{c}) \mathbf{s}^i,$$

and we will use  $\alpha_i = \sigma_i (\mathbf{r}^i)^\top \mathbf{c}$  to denote the coefficients of  $\mathbf{c}$  against the singular vectors of  $\mathbf{A}$ . Note that  $\sigma_i = 0$  implies  $\alpha_i = 0$ . The norm condition on  $\mathbf{A}^\top \mathbf{c}$  also means  $\sum_i \alpha_i^2 = \sum_{i:\sigma_i \neq 0} \alpha_i^2 \geq 1$ , which then gives

$$\begin{aligned}
\|\Pi_{\mathbf{A}} \mathbf{c}\|_2^2 &= \mathbf{c}^\top \mathbf{A} (\mathbf{A}^\top \mathbf{A})^\dagger \mathbf{A}^\top \mathbf{c} \\
&= \left( \sum_i \alpha_i (\mathbf{s}^i)^\top \right) \left( \sum_{i:\sigma_i \neq 0} \frac{1}{\sigma_i^2} \mathbf{s}^i (\mathbf{s}^i)^\top \right) \left( \sum_i \alpha_i \mathbf{s}^i \right) = \sum_{i:\sigma_i \neq 0} \frac{\alpha_i^2}{\sigma_i^2} \geq \frac{1}{\sigma_{\max}^2(\mathbf{A})}.
\end{aligned}$$

This completes the proof.  $\square$

*Proof of Claim 3.1.2.* The first inequality is trivial.

The second inequality follows the relations of norms,

$$\|\mathbf{c}\|_2 \leq \sqrt{n} \|\mathbf{c}\|_\infty \leq \sqrt{s}U.$$

For the third one, by the definition of Frobenius norm,

$$\|\mathbf{A}\|_F^2 = \sum_{i,j} \mathbf{A}_{ij}^2 = \sum_{i=1}^{\min\{m,n\}} \lambda_i(\mathbf{A}^\top \mathbf{A}).$$

Thus,

$$\lambda_{\max}(\mathbf{A}^\top \mathbf{A}) \leq \|\mathbf{A}\|_F^2 \leq sU^2,$$

and

$$\lambda_{\max}(\mathbf{A}^\top \mathbf{A}) \geq \frac{\|\mathbf{A}\|_F^2}{\min\{m,n\}} \geq \frac{1}{sU^2}.$$

By the definition of  $\kappa(\mathbf{A}^\top \mathbf{A})$ ,

$$\lambda_{\min}(\mathbf{A}^\top \mathbf{A}) = \frac{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})}{\kappa(\mathbf{A}^\top \mathbf{A})} \geq \frac{1}{sK^2U^2}.$$

This completes the proof. □

## REFERENCES

- [1] N. Megiddo, “On solving the linear programming problem approximately,” *Contemporary mathematics*, vol. 114, pp. 35–50, 1990.
- [2] J. Renegar, “Incorporating condition measures into the complexity theory of linear programming,” *SIAM Journal on Optimization*, vol. 5, no. 3, pp. 506–524, 1995.
- [3] N. Karmarkar, “A new polynomial-time algorithm for linear programming,” in *Proceedings of the sixteenth annual ACM symposium on Theory of computing*, ACM, 1984, pp. 302–311.
- [4] Y. T. Lee and A. Sidford, “Path finding methods for linear programming: solving linear programs in  $\tilde{O}(\sqrt{\text{rank}})$  iterations and faster algorithms for maximum flow,” in *Foundations of Computer Science (FOCS), 2014 IEEE 55th Annual Symposium on*, Available at <http://arxiv.org/abs/1312.6677> and <http://arxiv.org/abs/1312.6713>, IEEE, 2014, pp. 424–433.
- [5] L. G. Khachiyan, “A polynomial algorithm in linear programming,” in *Doklady Akademii Nauk SSSR*, vol. 244, 1979, pp. 1093–1096.
- [6] V. Klee and G. J. Minty, “How good is the simplex algorithm,” WASHINGTON UNIV SEATTLE DEPT OF MATHEMATICS, Tech. Rep., 1970.
- [7] D. Spielman and S.-H. Teng, “Smoothed analysis of algorithms: why the simplex algorithm usually takes polynomial time,” in *Proceedings of the thirty-third annual ACM symposium on Theory of computing*, ACM, 2001, pp. 296–305.
- [8] D. A. Spielman, “Nsf award 1562041: generalized algebraic graph theory: algorithms and analysis,” *ALGORITHMIC FOUNDATIONS*, 2016.
- [9] D. A. Spielman and S.-H. Teng, “Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems,” in *Proceedings of the Thirty-sixth Annual ACM Symposium on Theory of Computing*, ser. STOC ’04, Chicago, IL, USA: ACM, 2004, pp. 81–90, ISBN: 1-58113-852-0.
- [10] R. Kyng, Y. T. Lee, R. Peng, S. Sachdeva, and D. A. Spielman, “Sparsified cholesky and multigrid solvers for connection laplacians,” in *Proceedings of the Forty-eighth Annual ACM Symposium on Theory of Computing*, ser. STOC ’16, Cambridge, MA, USA: ACM, 2016, pp. 842–850, ISBN: 978-1-4503-4132-5.



- [11] M. B. Cohen, J. Kelner, J. Peebles, R. Peng, A. B. Rao, A. Sidford, and A. Vladu, “Almost-linear-time algorithms for markov chains and new spectral primitives for directed graphs,” in *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, ser. STOC 2017, Montreal, Canada: ACM, 2017, pp. 410–419, ISBN: 978-1-4503-4528-6.
- [12] M. B. Cohen, J. Kelner, R. Kyng, J. Peebles, R. Peng, A. B. Rao, and A. Sidford, “Solving directed laplacians in nearly linear time through sparse lu factorizations,” in *Foundations of Computer Science (FOCS), 2018 IEEE 50th Annual Symposium on*, IEEE, 2018.
- [13] N. E. Young, “Sequential and parallel algorithms for mixed packing and covering,” in *42nd Annual Symposium on Foundations of Computer Science, FOCS 2001, 14-17 October 2001, Las Vegas, Nevada, USA, 2001*, pp. 538–546.
- [14] V. Strassen, “Gaussian elimination is not optimal,” *Numerische mathematik*, vol. 13, no. 4, pp. 354–356, 1969.
- [15] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to algorithms*. MIT press, 2009.
- [16] V. V. Williams, “Multiplying matrices faster than coppersmith-winograd,” in *Proceedings of the Forty-fourth Annual ACM Symposium on Theory of Computing*, ser. STOC ’12, New York, New York, USA: ACM, 2012, pp. 887–898, ISBN: 978-1-4503-1245-5.
- [17] F. Le Gall, “Powers of tensors and fast matrix multiplication,” in *Proceedings of the 39th international symposium on symbolic and algebraic computation*, Available at: <https://arxiv.org/abs/1401.7714>, ACM, 2014, pp. 296–303.
- [18] R. Raz, “On the complexity of matrix product,” in *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*, ACM, 2002, pp. 144–151.
- [19] M. R. Hestenes and E. Stiefel, *Methods of conjugate gradients for solving linear systems*, 1. NBS Washington, DC, 1952, vol. 49.
- [20] C. Musco, C. Musco, and A. Sidford, “Stability of the lanczos method for matrix function approximation,” in *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, Society for Industrial and Applied Mathematics, 2018, pp. 1605–1624.
- [21] A. Madry, “Navigating central path with electrical flows: from flows to matchings, and back,” in *Proceedings of the 2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, ser. FOCS ’13, Washington, DC, USA: IEEE Computer Society, 2013, pp. 253–262, ISBN: 978-0-7695-5135-7.

- [22] —, “Computing maximum flow with augmenting electrical flows,” in *Proceedings of the 2016 IEEE 57th Annual Symposium on Foundations of Computer Science*, ser. FOCS ’16, Washington, DC, USA: IEEE Computer Society, 2016.
- [23] M. B. Cohen, A. Madry, P. Sankowski, and A. Vladu, “Negative-weight shortest paths and unit capacity minimum cost flow in  $\tilde{O}(M^{10/7} \log W)$  time: (extended abstract),” in *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, ser. SODA ’17, Barcelona, Spain: Society for Industrial and Applied Mathematics, 2017, pp. 752–771.
- [24] S. I. Daich and D. A. Spielman, “Faster approximate lossy generalized flow via interior point algorithms,” in *Proceedings of the 40th annual ACM symposium on Theory of computing*, ser. STOC ’08, Available at <http://arxiv.org/abs/0803.0988>, Victoria, British Columbia, Canada: ACM, 2008, pp. 451–460, ISBN: 978-1-60558-047-0.
- [25] S.-H. Teng, “The Laplacian Paradigm: Emerging Algorithms for Massive Graphs,” in *Theory and Applications of Models of Computation*, 2010, pp. 2–14.
- [26] I. Koutis, G. L. Miller, and R. Peng, “Approaching optimality for solving SDD linear systems,” in *Proceedings of the 2010 IEEE 51st Annual Symposium on Foundations of Computer Science*, ser. FOCS ’10, Available at <http://arxiv.org/abs/1003.2958>, Washington, DC, USA: IEEE Computer Society, 2010, pp. 235–244, ISBN: 978-0-7695-4244-7.
- [27] —, “A nearly- $m \log n$  time solver for SDD linear systems,” in *Proceedings of the 2011 IEEE 52nd Annual Symposium on Foundations of Computer Science*, ser. FOCS ’11, Available at <http://arxiv.org/abs/1102.4842>, Washington, DC, USA: IEEE Computer Society, 2011, pp. 590–598, ISBN: 978-0-7695-4571-4.
- [28] J. A. Kelner, L. Orecchia, A. Sidford, and Z. A. Zhu, “A simple, combinatorial algorithm for solving sdd systems in nearly-linear time,” in *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, ACM, 2013, pp. 911–920.
- [29] R. Peng and D. A. Spielman, “An efficient parallel solver for sdd linear systems,” in *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, ACM, 2014, pp. 333–342.
- [30] M. B. Cohen, R. Kyng, G. L. Miller, J. W. Pachocki, R. Peng, A. B. Rao, and S. C. Xu, “Solving sdd linear systems in nearly  $m \log 1/2 n$  time,” in *Proceedings of the 46th Annual ACM Symposium on Theory of Computing*, ACM, 2014, pp. 343–352.
- [31] R. Kyng and S. Sachdeva, “Approximate gaussian elimination for laplacians-fast, sparse, and simple,” in *Foundations of Computer Science (FOCS), 2016 IEEE 57th Annual Symposium on*, IEEE, 2016, pp. 573–582.

- [32] H. H. Chin, A. Madry, G. L. Miller, and R. Peng, “Runtime guarantees for regression problems,” in *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, ser. ITCS ’13, Berkeley, California, USA: ACM, 2013, pp. 269–282, ISBN: 978-1-4503-1859-4.
- [33] R. Szeliski, R. Zabih, D. Scharstein, O. Veksler, V. Kolmogorov, A. Agarwala, M. Tappen, and C. Rother, “A comparative study of energy minimization methods for markov random fields with smoothness-based priors,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, pp. 1068–1080, 2008.
- [34] Y. Wang, J. Yang, W. Yin, and Y. Zhang, “A new alternating minimization algorithm for total variation image reconstruction,” *SIAM Journal on Imaging Sciences*, vol. 1, no. 3, pp. 248–272, 2008.
- [35] M. Arie-Nachimson, S. Z. Kovalsky, I. Kemelmacher-Shlizerman, A. Singer, and R. Basri, “Global motion estimation from point matches,” in *3D Imaging, Modeling, Processing, Visualization and Transmission (3DIMPVT), 2012 Second International Conference on*, IEEE, 2012, pp. 81–88.
- [36] O. Ozyesil, A. Singer, and R. Basri, “Stable camera motion estimation using convex programming,” *SIAM Journal on Imaging Sciences*, vol. 8, no. 2, pp. 1220–1262, 2015.
- [37] M. B. Cohen, Y. T. Lee, G. Miller, J. Pachocki, and A. Sidford, “Geometric median in nearly linear time,” in *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing*, Available at: <https://arxiv.org/abs/1606.05225>, ACM, 2016, pp. 9–21.
- [38] A. Singer and Y. Shkolnisky, “Three-dimensional structure determination from common lines in cryo-EM by eigenvectors and semidefinite programming,” *SIAM journal on imaging sciences*, vol. 4, no. 2, pp. 543–572, 2011.
- [39] Y. Shkolnisky and A. Singer, “Viewing direction estimation in cryo-EM using synchronization,” *SIAM journal on imaging sciences*, vol. 5, no. 3, pp. 1088–1110, 2012.
- [40] Z. Zhao and A. Singer, “Rotationally invariant image representation for viewing direction classification in cryo-EM,” *Journal of structural biology*, vol. 186, no. 1, pp. 153–166, 2014.
- [41] B. Alexeev, A. S. Bandeira, M. Fickus, and D. G. Mixon, “Phase retrieval with polarization,” *SIAM Journal on Imaging Sciences*, vol. 7, no. 1, pp. 35–66, 2014.
- [42] S. Marchesini, Y.-C. Tu, and H.-t. Wu, “Alternating projection, ptychographic imaging and phase synchronization,” *arXiv preprint arXiv:1402.0550*, 2014.

- [43] G. Shklarski and S. Toledo, “Rigidity in finite-element matrices: sufficient conditions for the rigidity of structures and substructures,” *SIAM Journal on Matrix Analysis and Applications*, vol. 30, no. 1, pp. 7–40, 2008.
- [44] S. I. Daitch and D. A. Spielman, “Support-graph preconditioners for 2-dimensional trusses,” *arXiv preprint cs/0703119*, 2007.
- [45] A. Singer and H.-T. Wu, “Vector diffusion maps and the connection laplacian,” *Communications on pure and applied mathematics*, vol. 65, no. 8, pp. 1067–1144, 2012.
- [46] A. S. Bandeira, A. Singer, and D. A. Spielman, “A cheeger inequality for the graph connection laplacian,” *SIAM Journal on Matrix Analysis and Applications*, vol. 34, no. 4, pp. 1611–1630, 2013.
- [47] A. George, “Nested dissection of a regular finite element mesh,” *SIAM Journal on Numerical Analysis*, vol. 10, no. 2, pp. 345–363, 1973.
- [48] R. J. Lipton, D. J. Rose, and R. E. Tarjan, “Generalized nested dissection,” *SIAM journal on numerical analysis*, vol. 16, no. 2, pp. 346–358, 1979.
- [49] G. L. Miller and W. Thurston, “Separators in two and three dimensions,” in *Proceedings of the twenty-second annual ACM symposium on Theory of computing*, ACM, 1990, pp. 300–309.
- [50] N. Alon and R. Yuster, “Solving linear systems through nested dissection,” in *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, IEEE, 2010, pp. 225–234.
- [51] M. Luby and N. Nisan, “A parallel approximation algorithm for positive linear programming,” in *Proceedings of the Twenty-Fifth Annual ACM Symposium on Theory of Computing, May 16-18, 1993, San Diego, CA, USA*, 1993, pp. 448–457.
- [52] S. Arora, E. Hazan, and S. Kale, “The multiplicative weights update method: a meta-algorithm and applications,” *Theory of Computing*, vol. 8, no. 1, pp. 121–164, 2012.
- [53] Z. Allen-Zhu and L. Orecchia, “Nearly-linear time positive LP solver with faster convergence rate,” in *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing*, ser. STOC ’15, Newer version available at <http://arxiv.org/abs/1411.1124>, 2015, pp. 229–236.
- [54] —, “Using optimization to break the epsilon barrier: a faster and simpler width-independent algorithm for solving positive linear programs in parallel,” in *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms*, ser. SODA ’15, 2015, pp. 1439–1456.

- [55] D. Wang, S. Rao, and M. W. Mahoney, “Unified acceleration method for packing and covering problems via diameter reduction,” in *43rd International Colloquium on Automata, Languages, and Programming, ICALP 2016, July 11-15, 2016, Rome, Italy*, 2016, 50:1–50:13.
- [56] R. Kyng and P. Zhang, “Hardness results for structured linear systems,” in *Foundations of Computer Science (FOCS), 2017 IEEE 58th Annual Symposium on*, IEEE, 2017, pp. 684–695.
- [57] R. Kyng, R. Peng, D. Wang, and P. Zhang, “Packing lps are hard to solve accurately, assuming linear equations are hard,” 2018, in submission.
- [58] L. Trevisan and F. Xhafa, “The parallel complexity of positive linear programming,” *Parallel Processing Letters*, vol. 8, no. 04, pp. 527–533, 1998.
- [59] R. Kyng, R. Peng, R. Schwieterman, and P. Zhang, “Incomplete nested dissection,” in *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, ACM, 2018, pp. 404–417.
- [60] M. W. Mahoney, S. Rao, D. Wang, and P. Zhang, “Approximating the solution to mixed packing and covering lps in parallel  $\tilde{O}(\epsilon^{-3})$  time,” in *LIPICs-Leibniz International Proceedings in Informatics*, Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, vol. 55, 2016.
- [61] D. Wang, M. W. Mahoney, N. Mohan, and S. Rao, “Faster parallel solver for positive linear programs via dynamically-bucketed selective coordinate descent,” *CoRR*, vol. abs/1511.06468, 2015.
- [62] L. N. Trefethen and D. Bau III, *Numerical linear algebra*. Siam, 1997, vol. 50.
- [63] K. L. Clarkson and D. P. Woodruff, “Low-rank approximation and regression in input sparsity time,” *Journal of the ACM (JACM)*, vol. 63, no. 6, p. 54, 2017.
- [64] E. G. Boman, D. Chen, O. Parekh, and S. Toledo, “On factor width and symmetric h-matrices,” *Linear algebra and its applications*, vol. 405, pp. 239–248, 2005.
- [65] D. Goldfarb and W. Yin, “Second-order cone programming methods for total variation-based image restoration,” *SIAM Journal on Scientific Computing*, vol. 27, no. 2, pp. 622–645, 2005.
- [66] M. B. Cohen, B. T. Fasy, G. L. Miller, A. Nayyeri, R. Peng, and N. Walkington, “Solving 1-laplacians in nearly linear time: collapsing and expanding a topological ball,” in *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2014, Portland, Oregon, USA, January 5-7, 2014*, Available

at: <https://www.cs.cmu.edu/~glmiller/Publications/Papers/CoFMNPW14.pdf>, 2014, pp. 204–216.

- [67] G. Barequet and S. Har-Peled, “Efficiently approximating the minimum-volume bounding box of a point set in three dimensions,” *J. Algorithms*, vol. 38, no. 1, pp. 91–109, Jan. 2001.
- [68] O. Axelsson, “A survey of preconditioned iterative methods for linear systems of algebraic equations,” *BIT Numerical Mathematics*, vol. 25, no. 1, pp. 165–187, 1985.
- [69] D. A. Spielman, “Spectral graph theory and its applications,” in *Foundations of Computer Science, 2007. FOCS’07. 48th Annual IEEE Symposium on*, IEEE, 2007, pp. 29–38.
- [70] S. J. Wright, *Primal-dual interior-point methods*. SIAM, 1997.
- [71] Y. Ye, *Interior point algorithms: theory and analysis*. John Wiley & Sons, 2011, vol. 44, Available at: <http://web.stanford.edu/~yye/main.ps>.
- [72] A. Nemirovski, “Interior point polynomial time methods in convex programming,” *Lecture notes*, 2004.
- [73] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004, Available at [https://web.stanford.edu/~boyd/cvxbook/bv\\_cvxbook.pdf](https://web.stanford.edu/~boyd/cvxbook/bv_cvxbook.pdf).
- [74] Y. T. Lee and A. Sidford, “Efficient inverse maintenance and faster algorithms for linear programming,” in *Foundations of Computer Science (FOCS), 2015 IEEE 56th Annual Symposium on*, Available at: <https://arxiv.org/abs/1503.01752>, IEEE, 2015, pp. 230–249.
- [75] A. Madry, “Faster approximation schemes for fractional multicommodity flow problems via dynamic graph algorithms,” in *STOC ’10: Proceedings of the 42nd ACM symposium on Theory of computing*, Cambridge, Massachusetts, USA: ACM, 2010, pp. 121–130, ISBN: 978-1-4503-0050-6.
- [76] D. Goldfarb and W. Yin, “Second-order cone programming methods for total variation-based image restoration,” *SIAM J. Sci. Comput.*, vol. 27, pp. 622–645, 2004.
- [77] G. N. Frederickson, “Fast algorithms for shortest paths in planar graphs, with applications,” *SIAM J. Comput.*, vol. 16, no. 6, pp. 1004–1022, 1987.
- [78] N. E. Young, “Nearly linear-time approximation schemes for mixed packing/covering and facility-location linear programs,” *CoRR*, vol. abs/1407.3015, 2014.